

What have we learned from deep representations for action recognition?

work with

Axel Pinz
Graz University of Technology



Richard P. Wildes
York University, Toronto



Andrew Zisserman
University of Oxford



Outline

- Two-Stream Architectures for Action Recognition

What have we learned in:

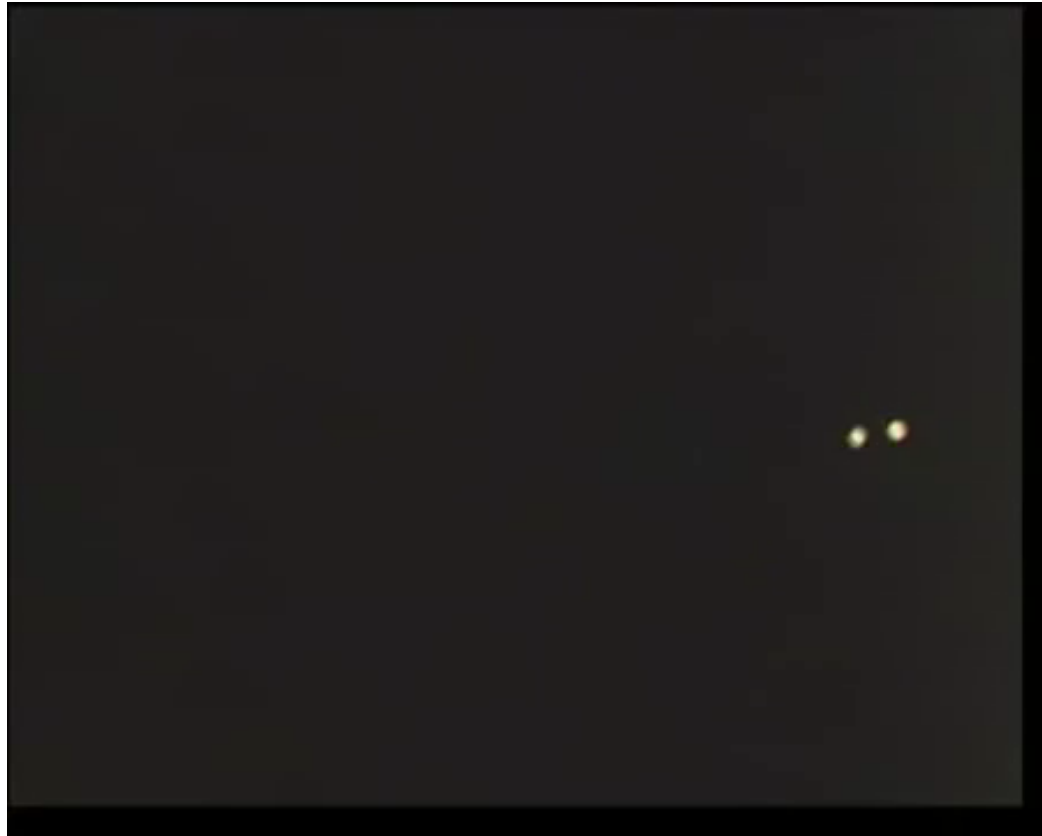
- Fusion of appearance and motion streams
- Long-term feature aggregation

- Visualization of Two-Stream representations

Intuitions for why:

- Explicit motion models perform better
- Fusion leads to good feature abstractions

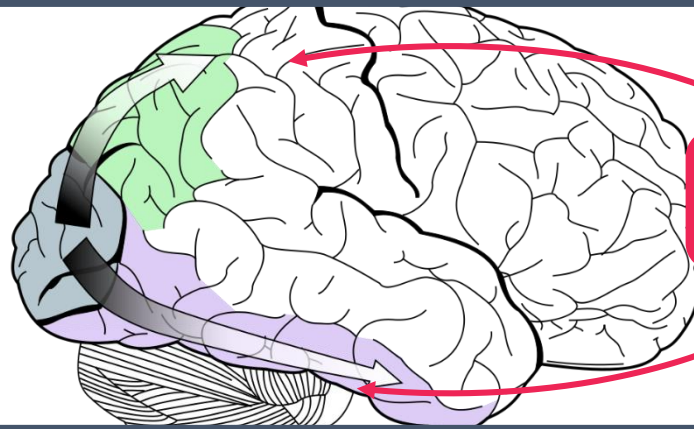
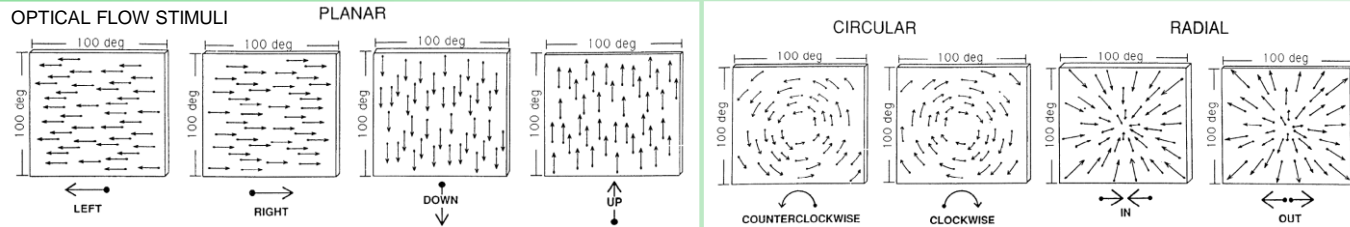
→ Amazing what the brain can do
without appearance information



Sources: Johansson, G. "Visual perception of biological motion and a model for its analysis." *Perception & Psychophysics*. 14(2):201-211. 1973.

Motivation: Separate visual pathways in nature

→ Dorsal stream ('where/how') recognizes motion and locates objects



→ “Interconnection”
e.g. in STS area

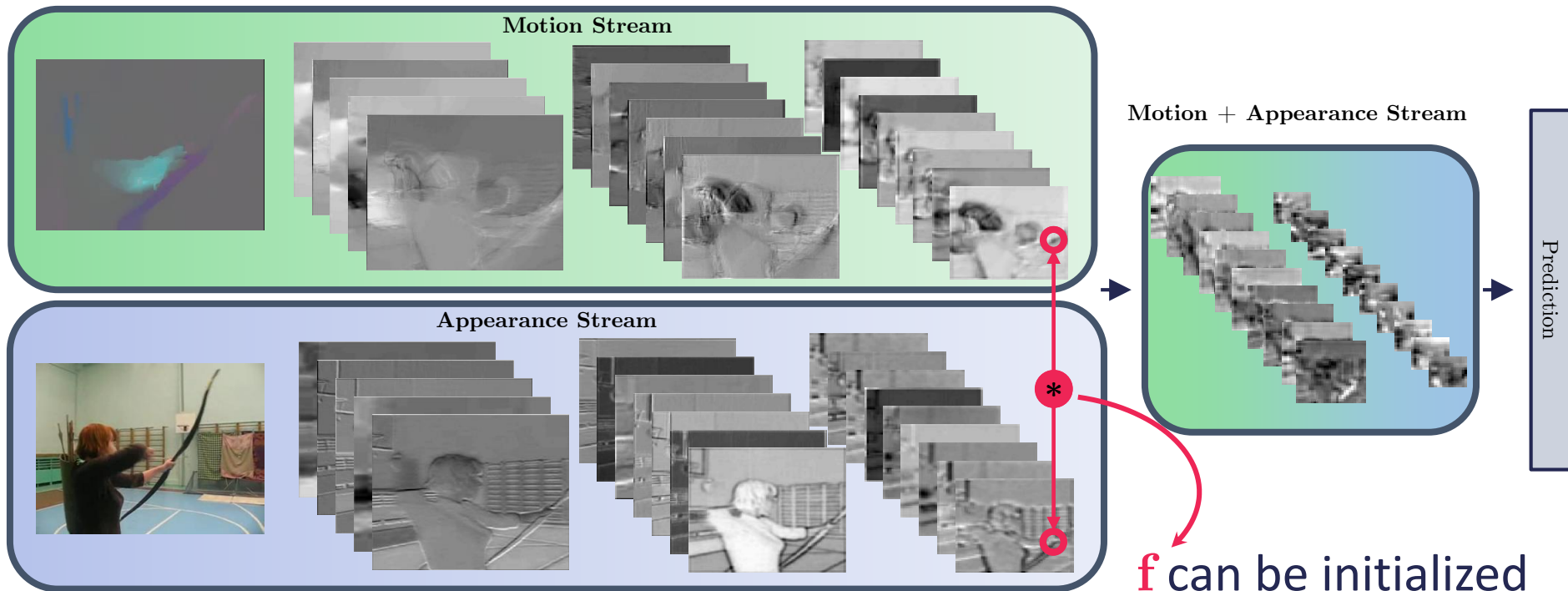
→ Ventral ('what') stream performs object recognition



Convolutional Two-Stream Network Fusion

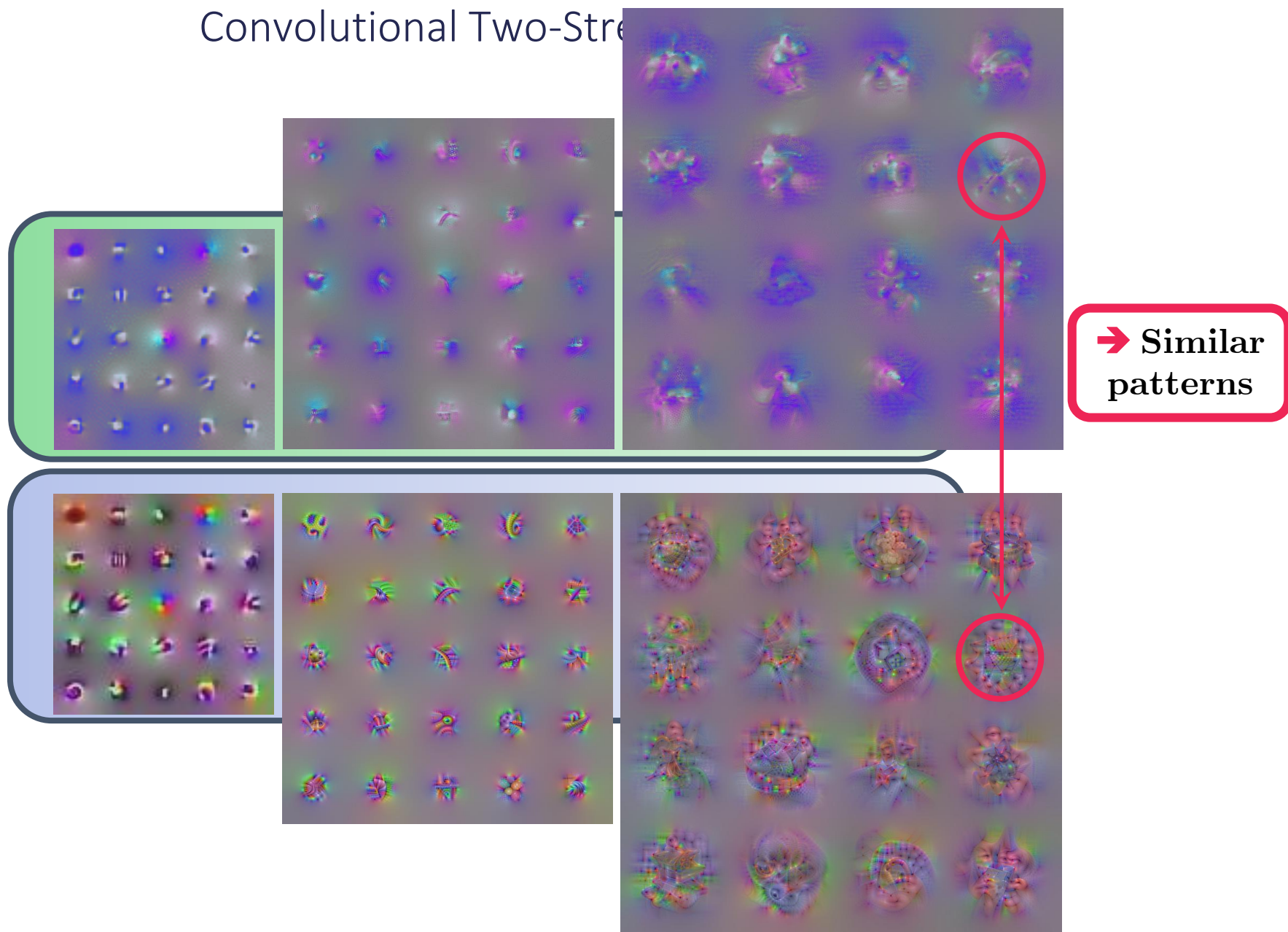
- We study a number of ways of **fusing** two-stream ConvNets

[Simonyan & Zisserman, NIPS'14]

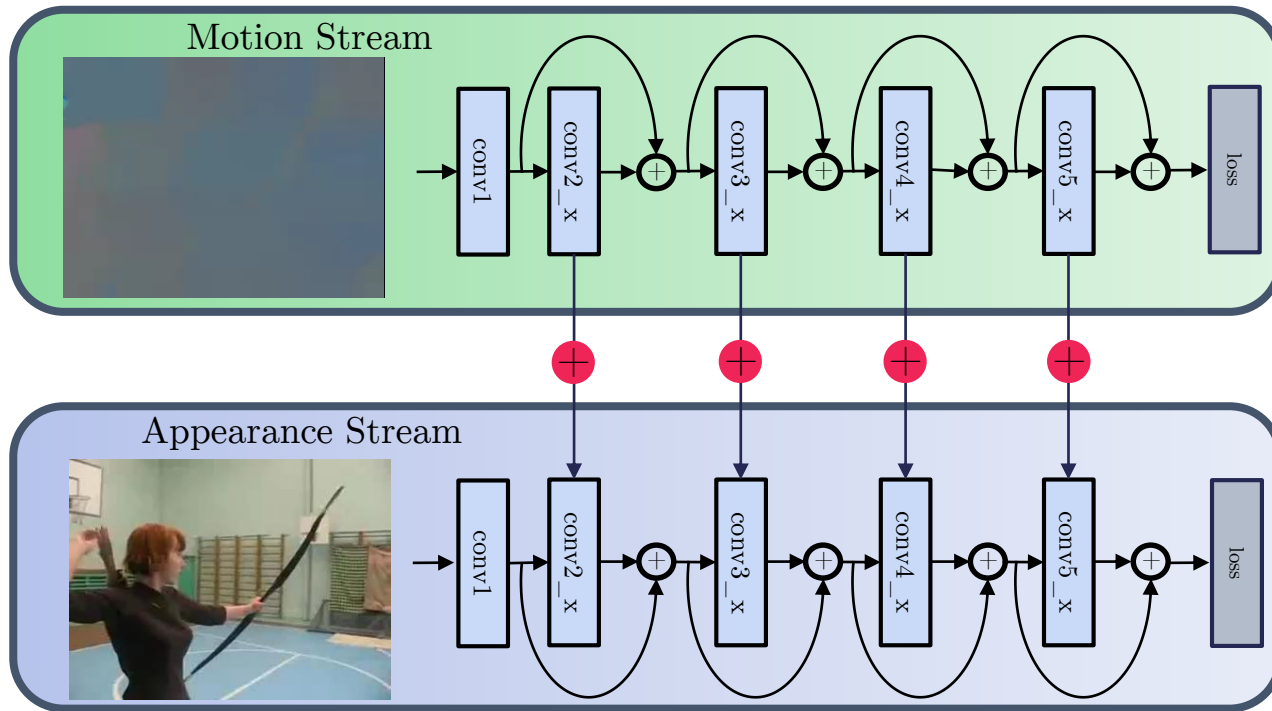


- Sum fusion works surprisingly well

Convolutional Two-Stream

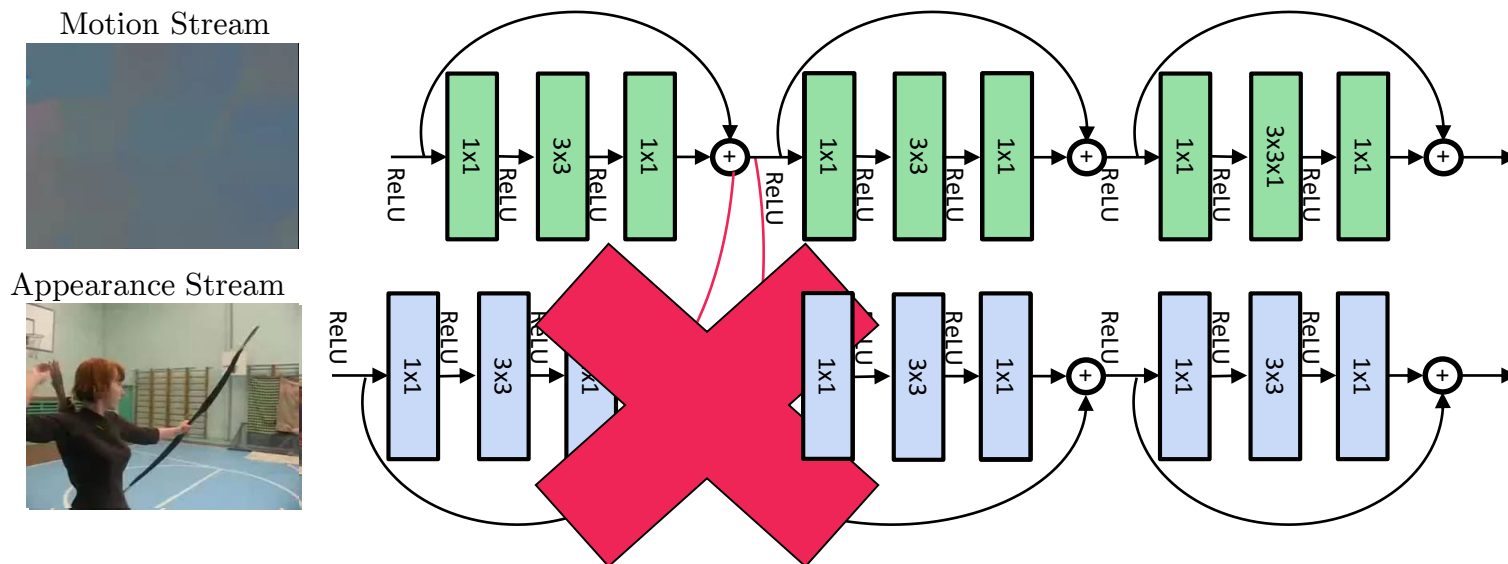


Spatiotemporal Residual Networks



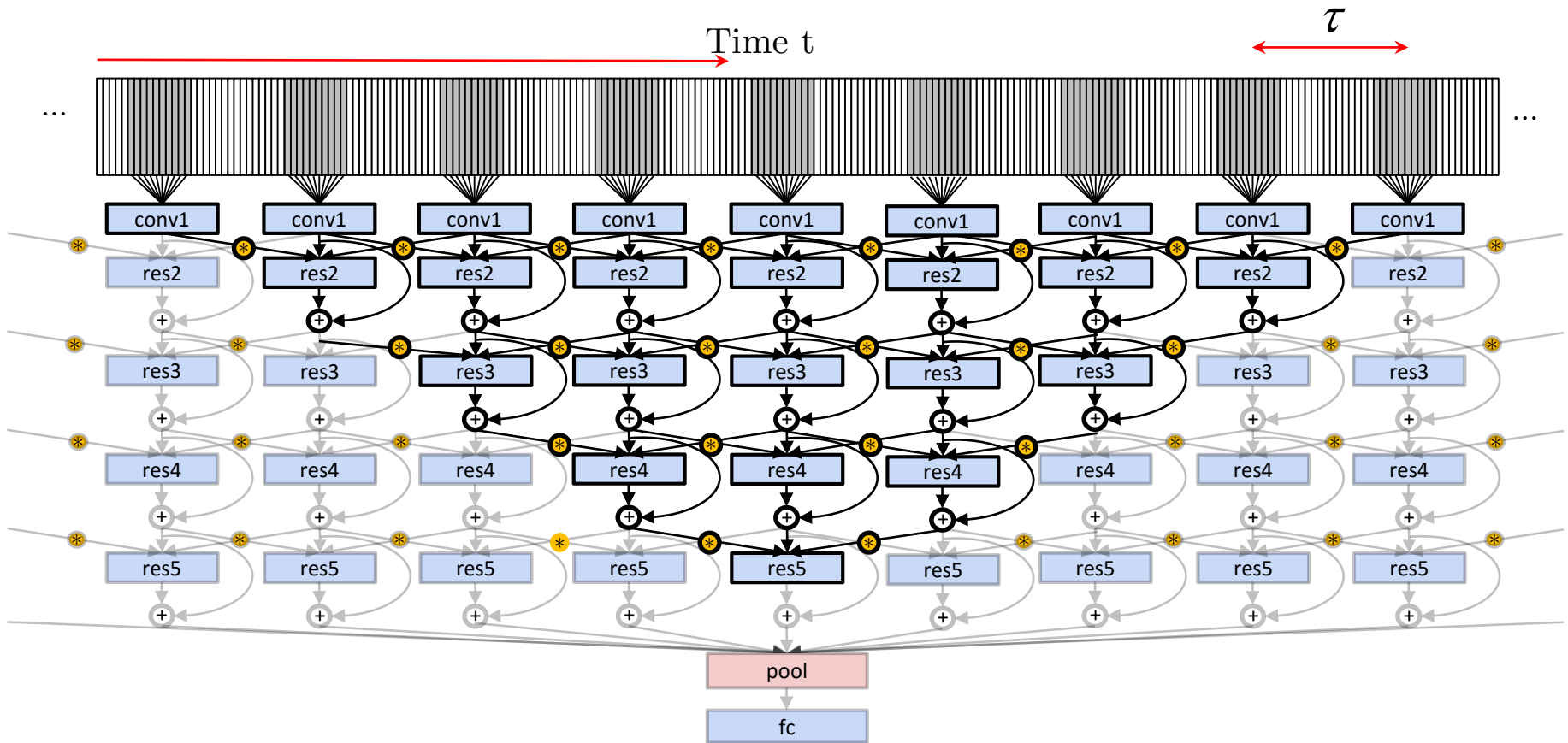
- ST-ResNet allows the hierarchical learning of spacetime features by **connecting** the appearance and motion channels of a two-stream architecture.
- Though, naive fusion does not work.

Fusing Two-Stream ResNets & Injecting Temporal Filters



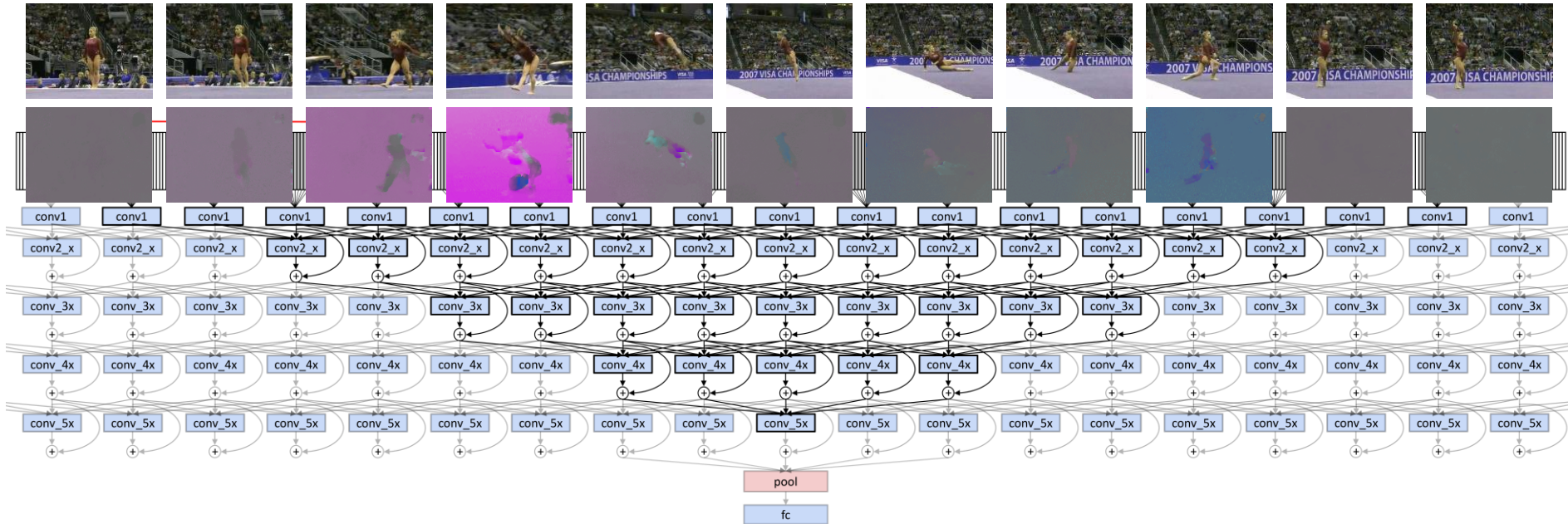
- ResNets for the spatiotemporal domain by introducing residual connections in two ways
 1. **Residuals** between the motion and appearance pathways to allow spatiotemporal interaction between the streams
 2. **Transformation** * of pretrained image ConvNets by filters initialized as residuals in time
- Our **most recent work** (@CVPR'17) reconsiders the combination these approaches more thoroughly to increase our understanding of how these techniques interact.

Increasing the temporal receptive field of ResNets



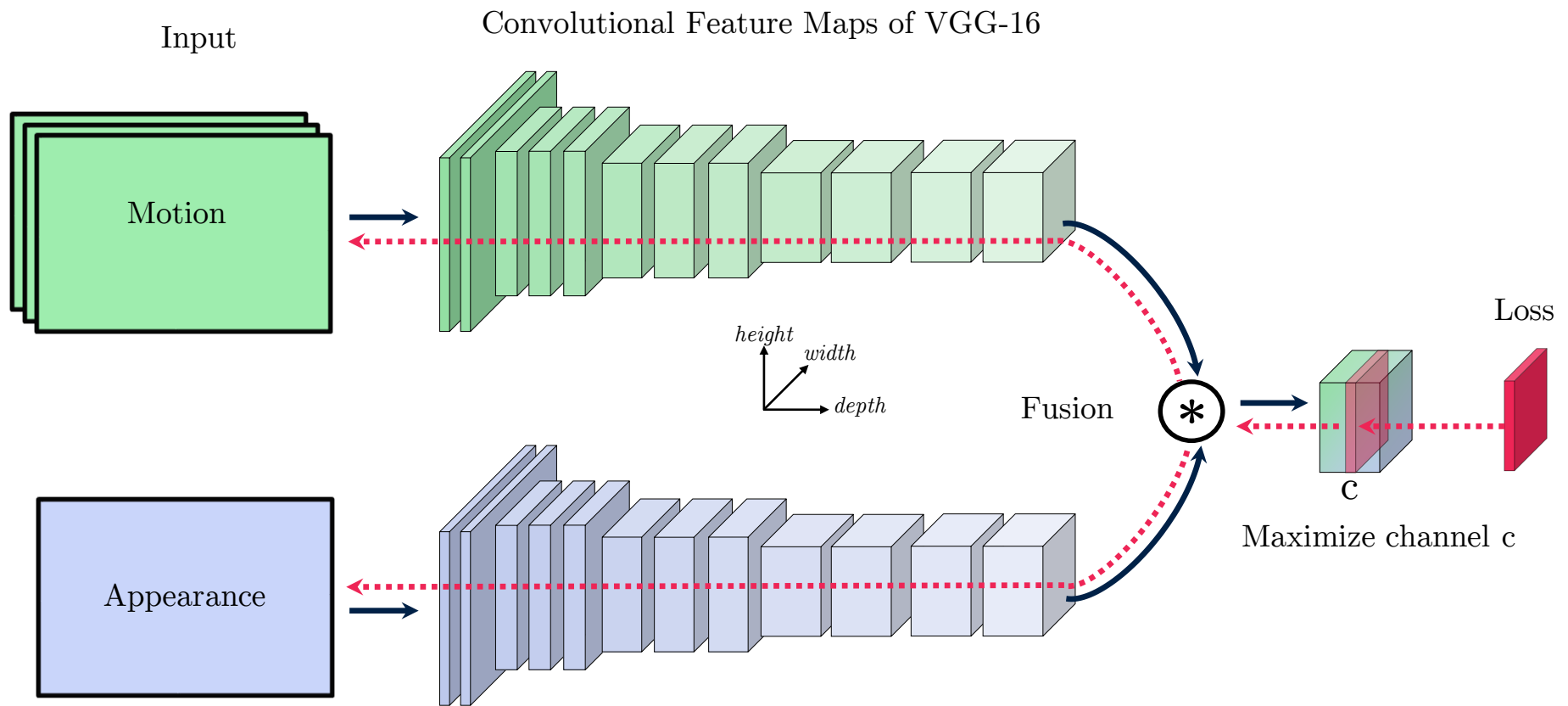
- The temporal receptive field is modulated by the temporal filters * and input stride τ

Transforming spatial filters to spatiotemporal ones

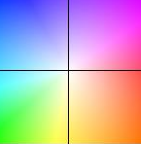


- Chaining temporal filters supports *hierarchical* learning of long-term correspondences between features of the appearance and motion stream.
- For example, if the stride is set to $\tau = 15$ frames and we transform 8 filters, a unit at conv5_3 sees a window of $17 \times 15 = 255$ frames.

Visualizing the learned representation



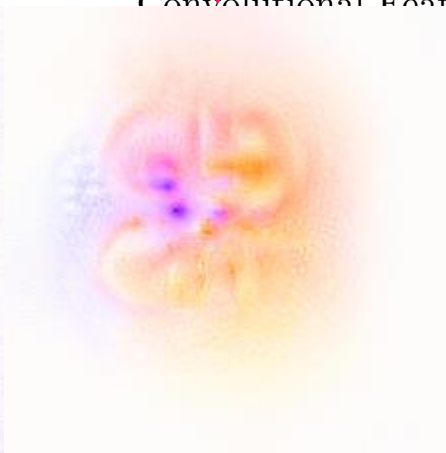
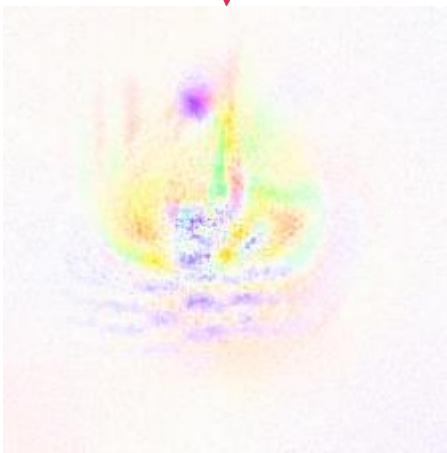
[Erhan et al. 2009]
[Simonyan et al. 2013]
[Mahendran & Vedaldi 2014]
[Yosinski et al. 2014]



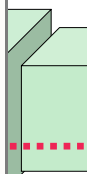
Visualizing the learned representation

Slow motion

Fast motion



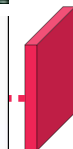
Convolutional Feature Map



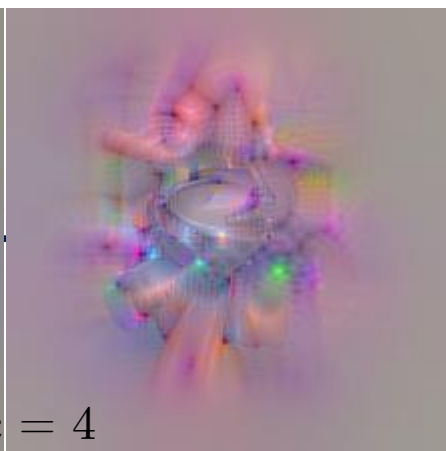
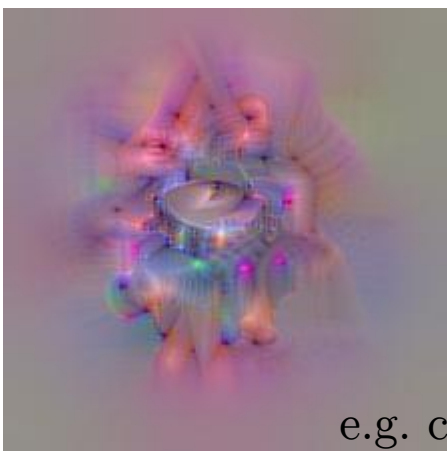
height
width



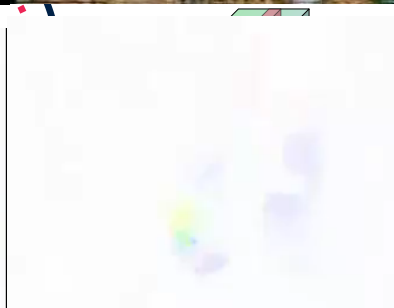
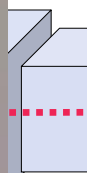
Loss



c



e.g. $c = 4$



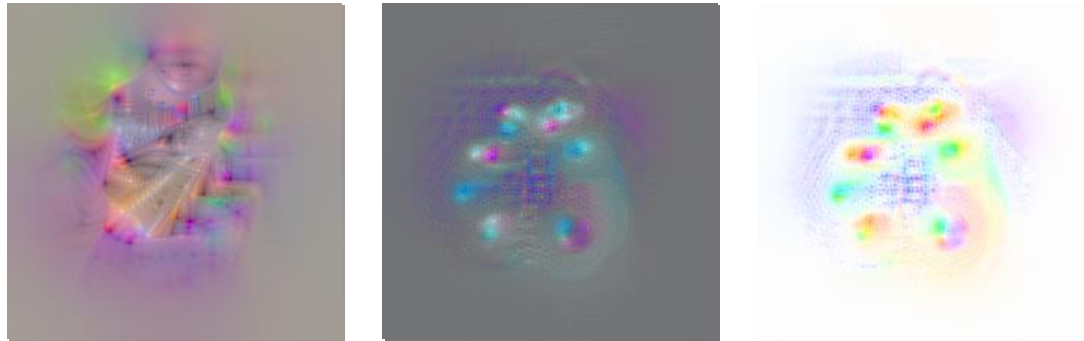
2014]



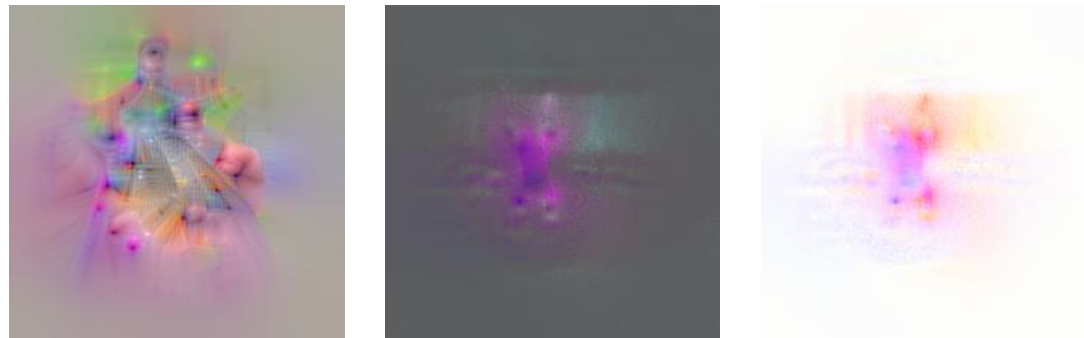
Filter #21 at conv5 fusion – a local Billiard neuron ?

→ (mostly) linear motion

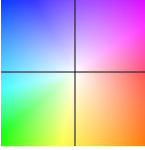
Slow motion



Fast motion

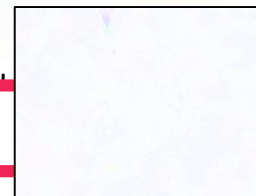
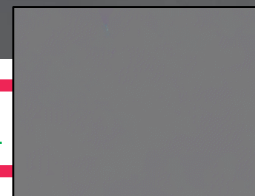
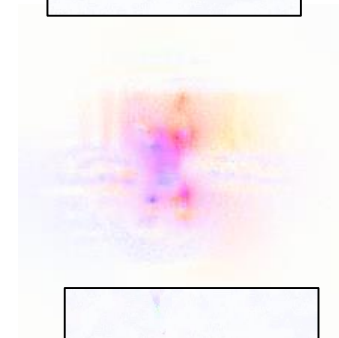
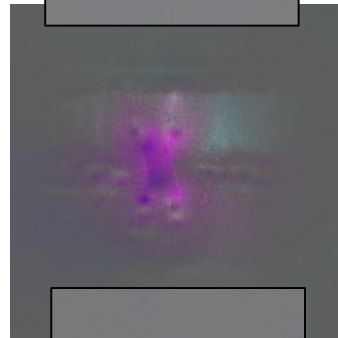
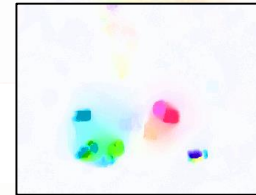
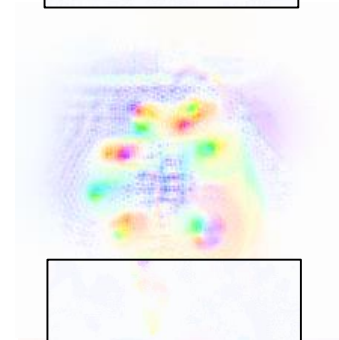
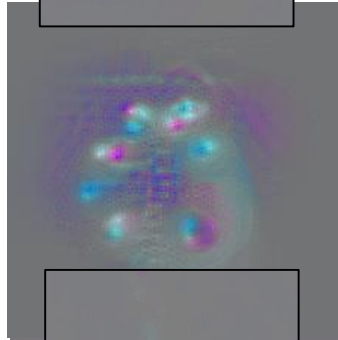
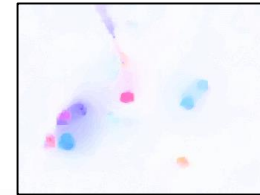
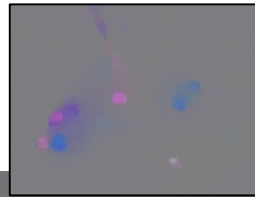


→ accelerating motion in all directions



Filter #21 at conv5 fusion – a local Billiard neuron ?

→ (mostly) linear



→ Test set
example
snippets

→ acceleration in all

↑
slower



Last layer

→ “Billiards”



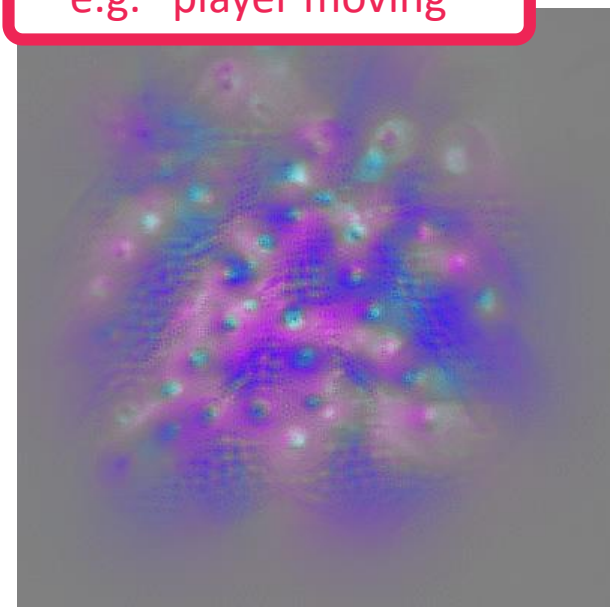
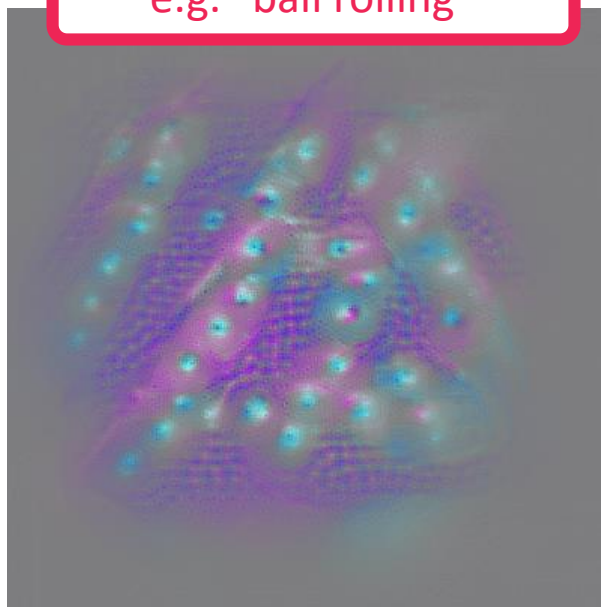
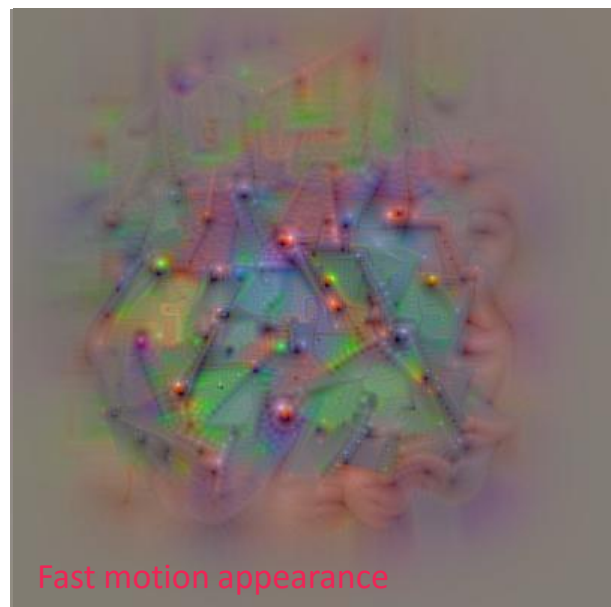
Appearance

Slow motion

e.g. “ball rolling”

Fast motion

e.g. “player moving”



sis 2017

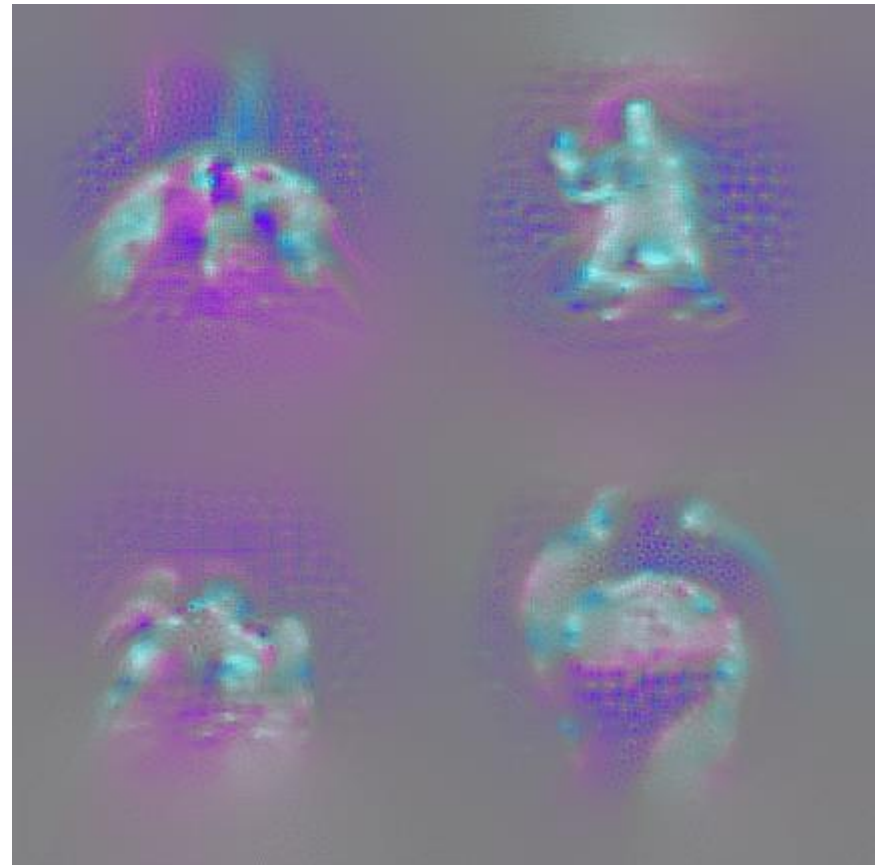
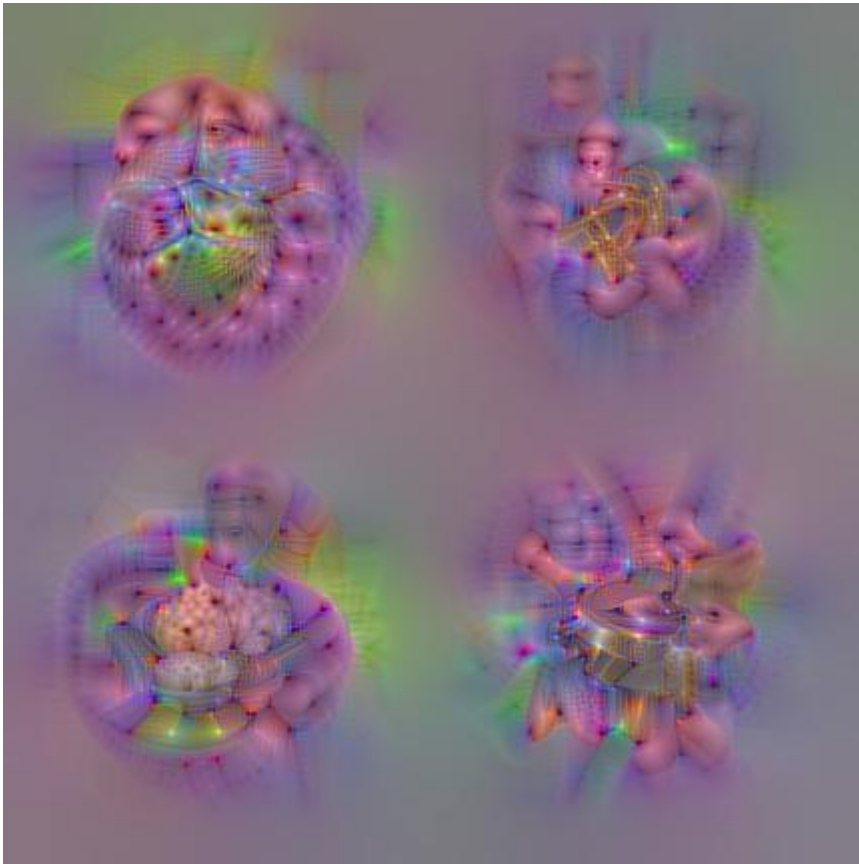


Going through the conv layers of VGG-16 (first four filters of each layer are shown)

Appearance

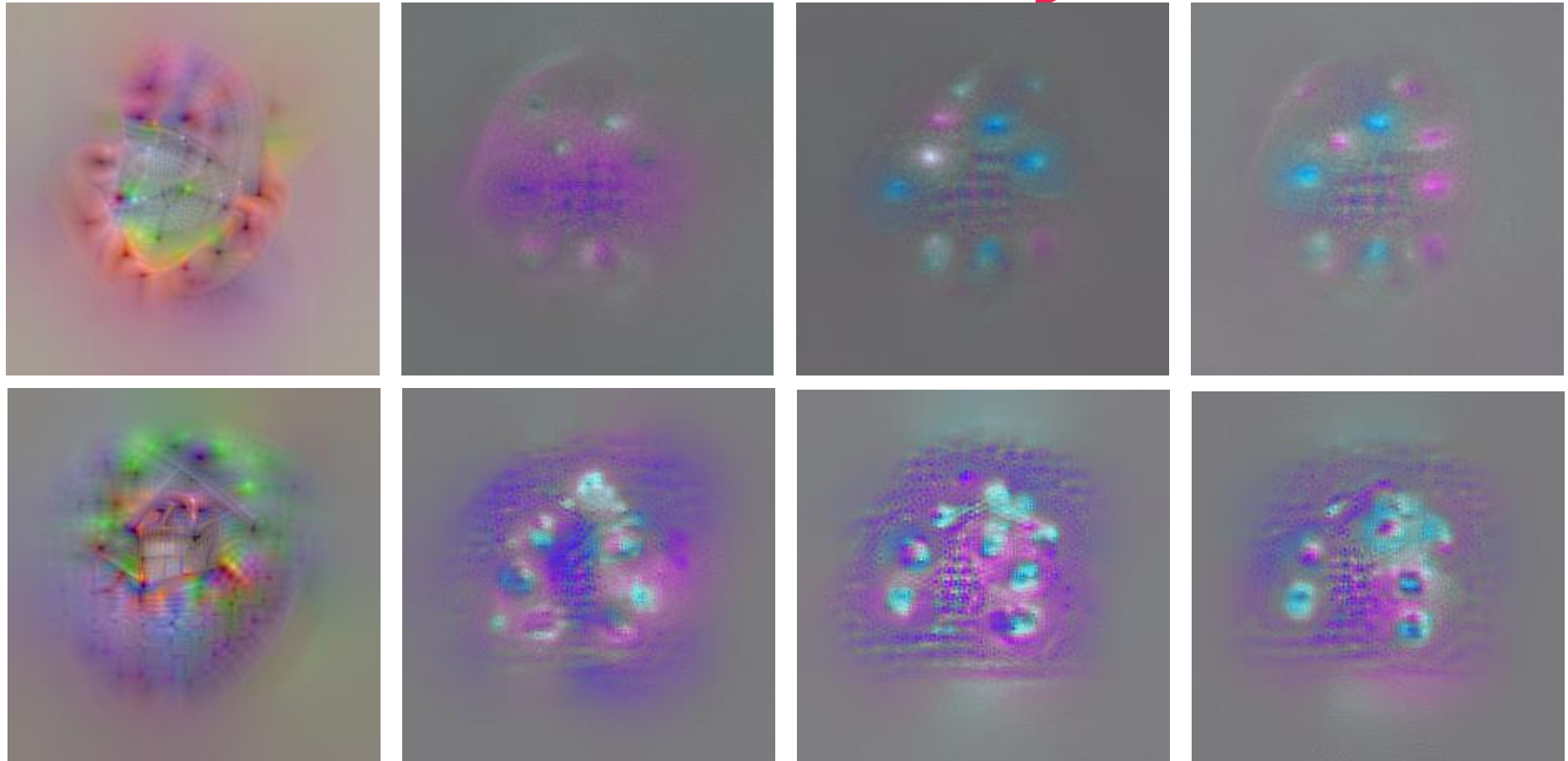
conv4_3 f1-4

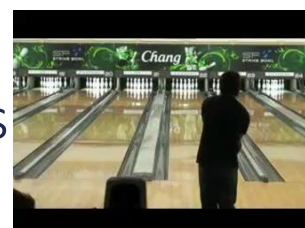
Slow motion



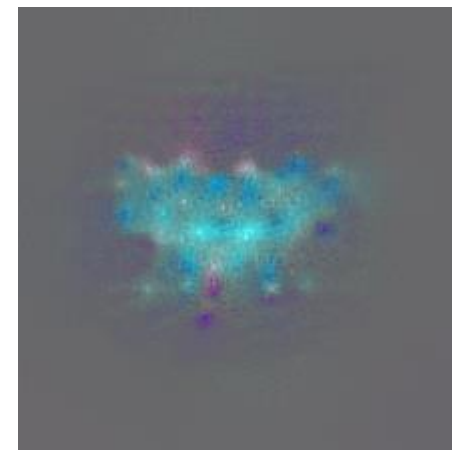
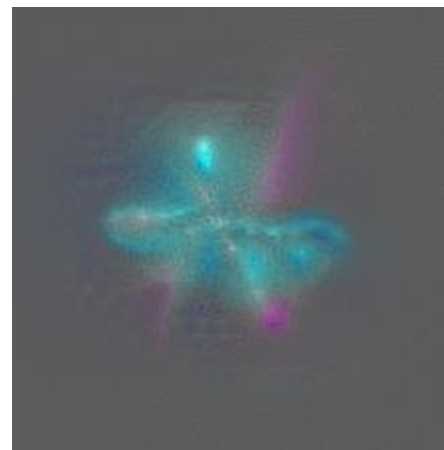
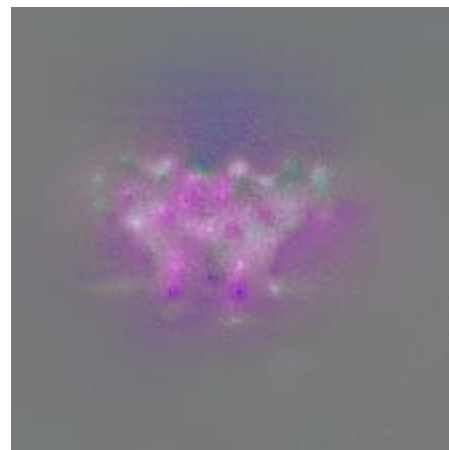
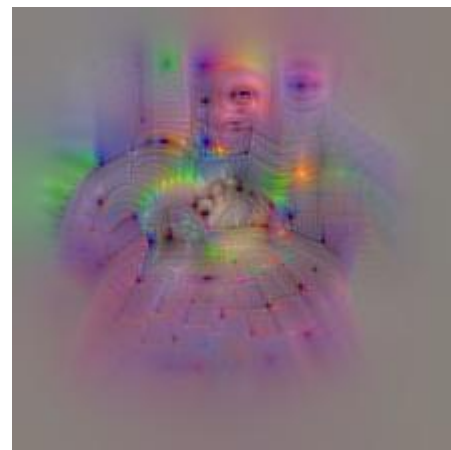
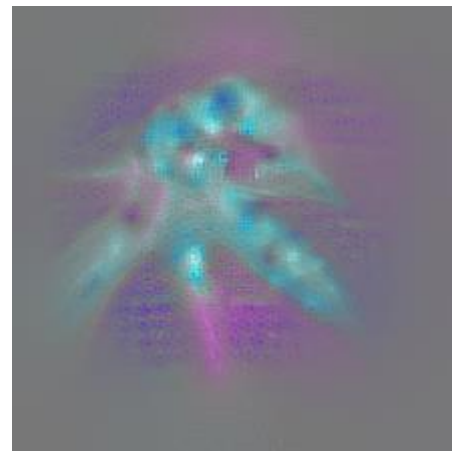
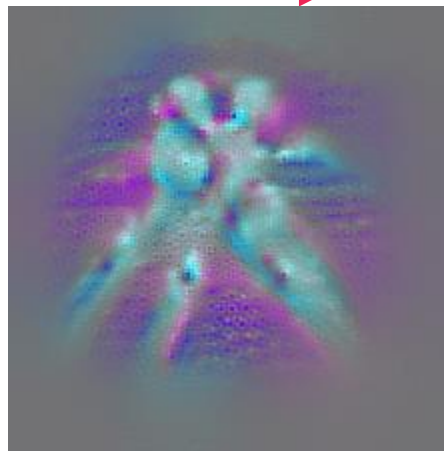
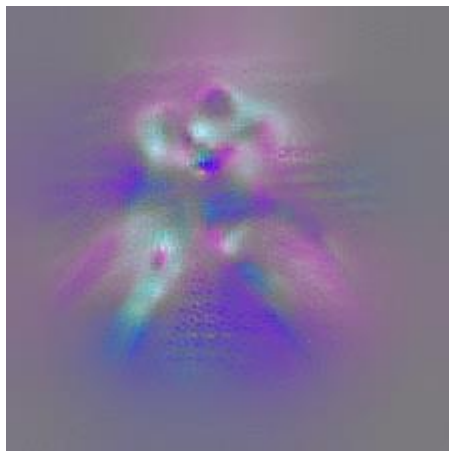


slower





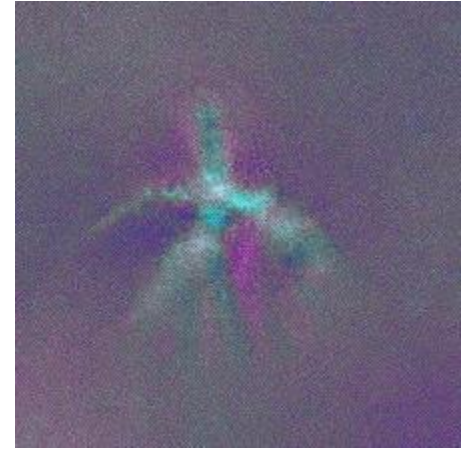
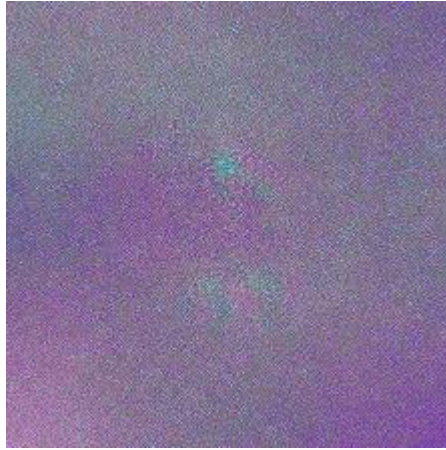
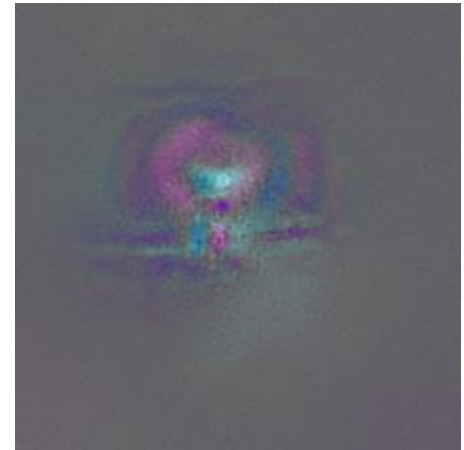
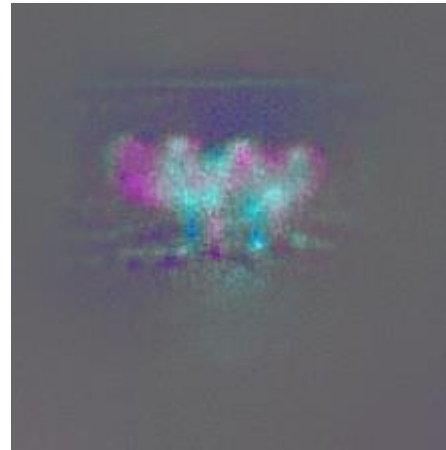
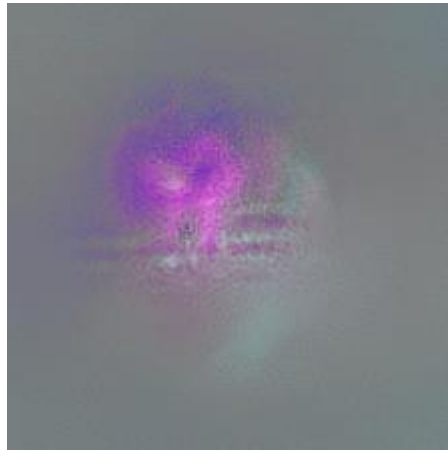
slower

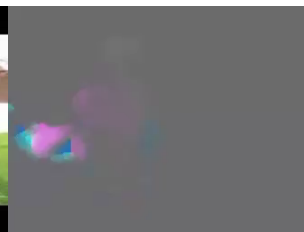
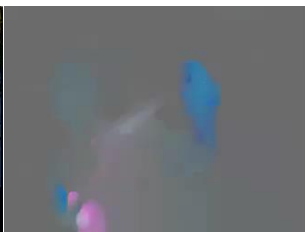




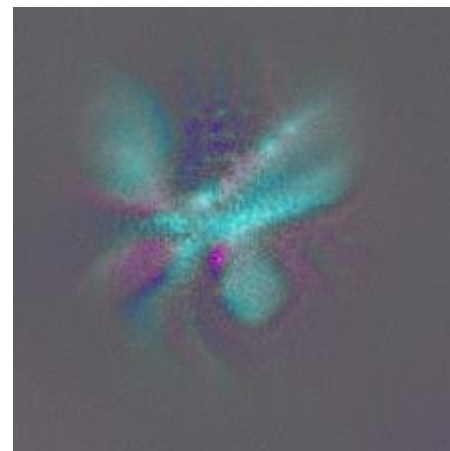
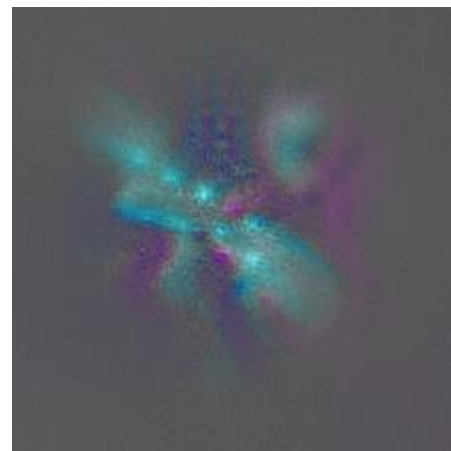
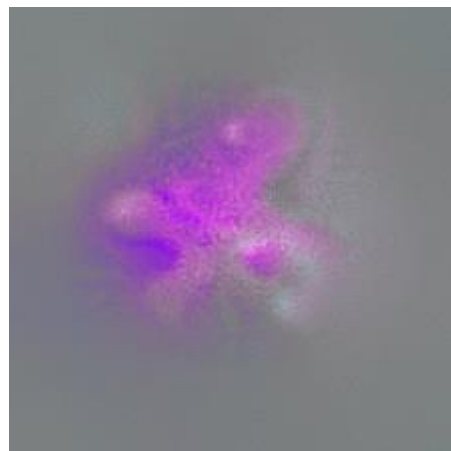
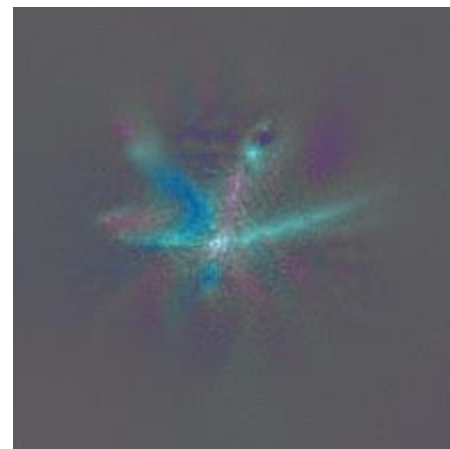
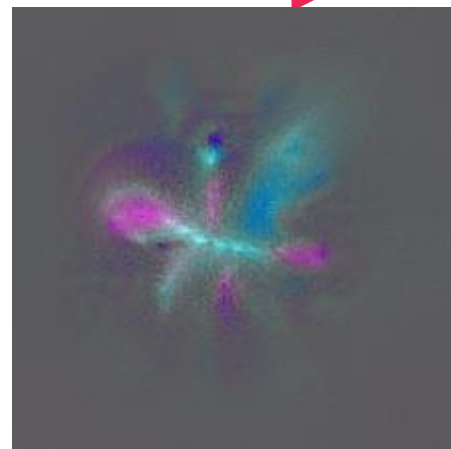
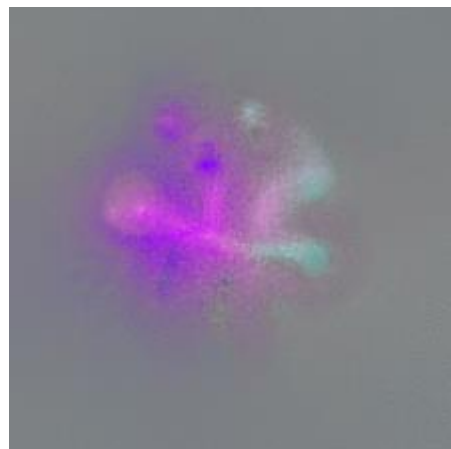
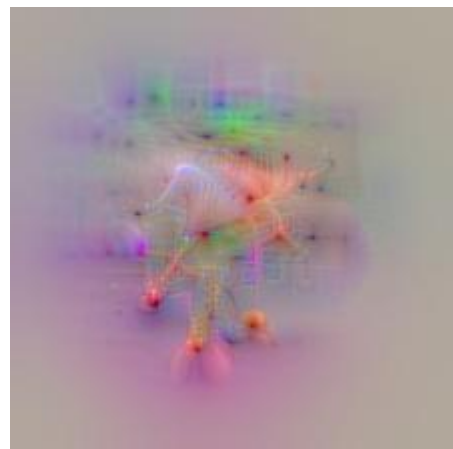
s a

slower



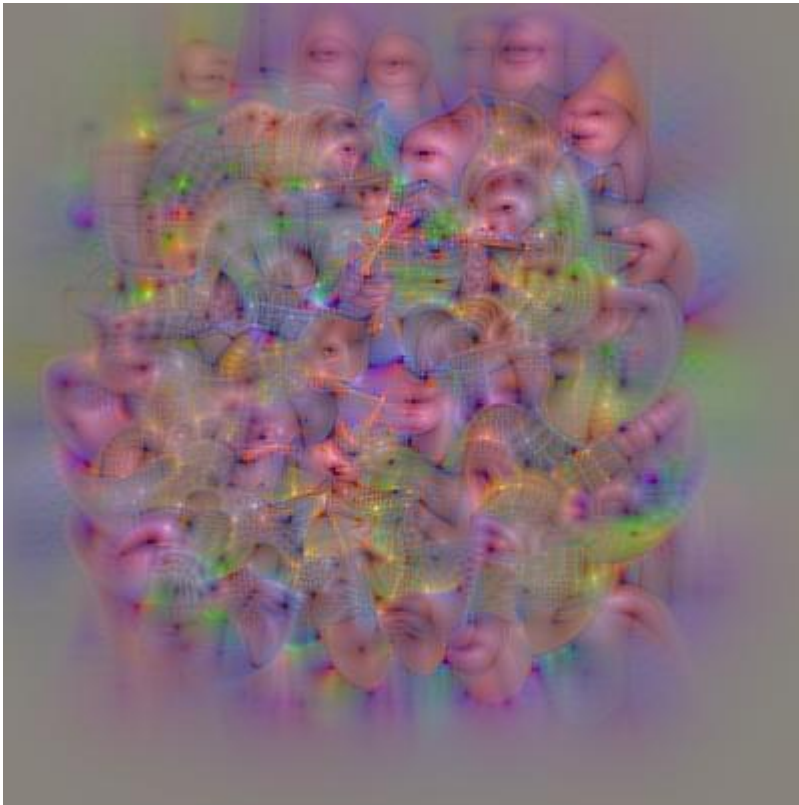


slower

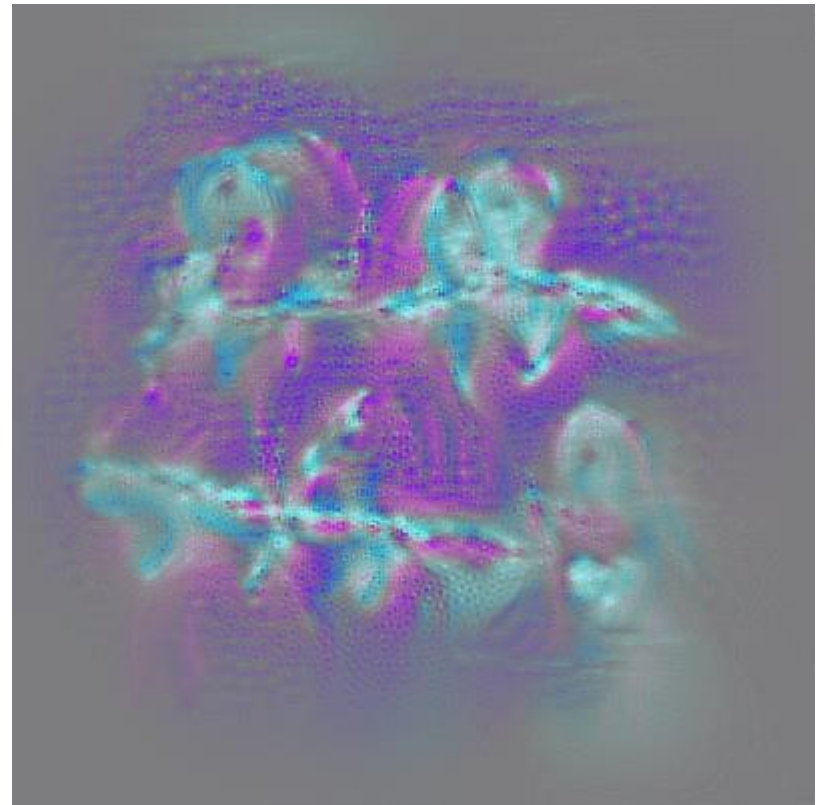


FC 6 (4096 features; RF 404x404)

Appearance

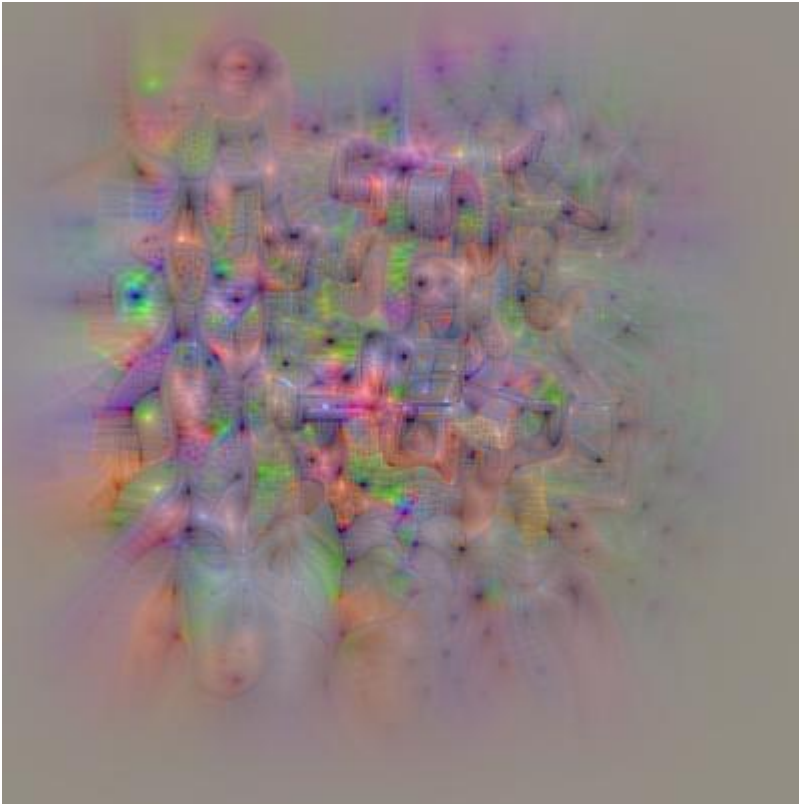


Slow motion

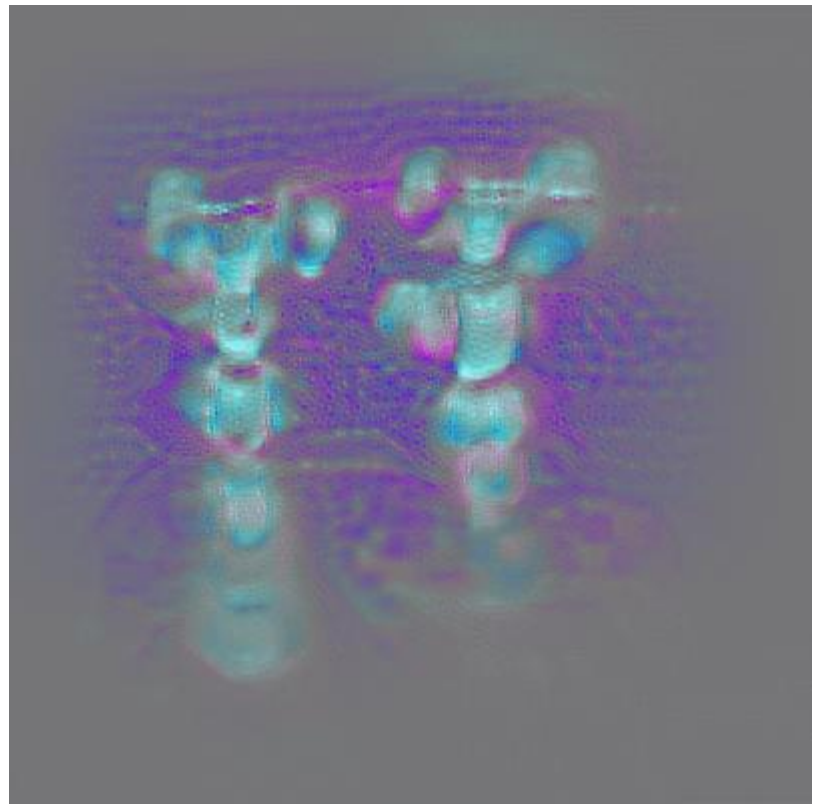


FC 7 (4096 features; RF 404x404)

Appearance



Slow motion

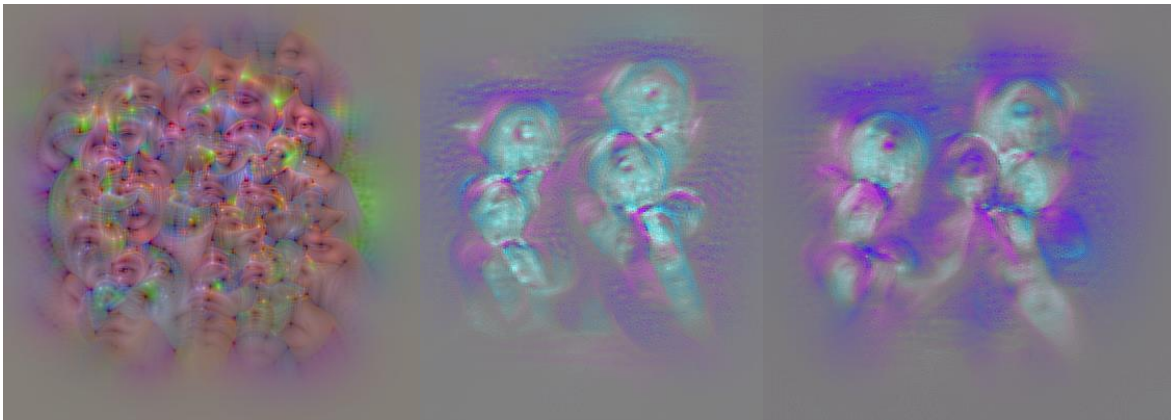


Explaining failure cases:

Appearance

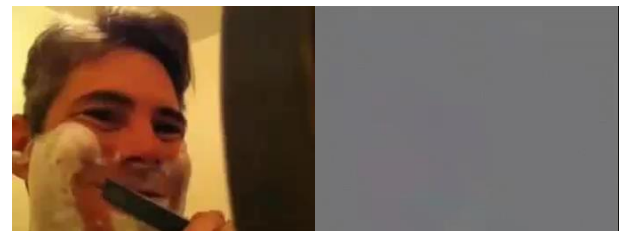
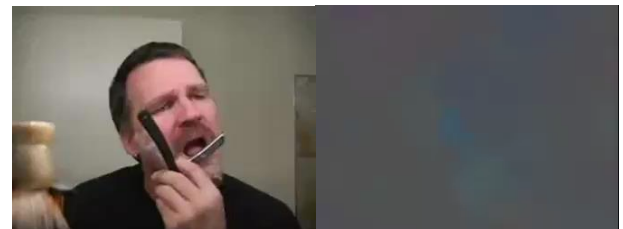
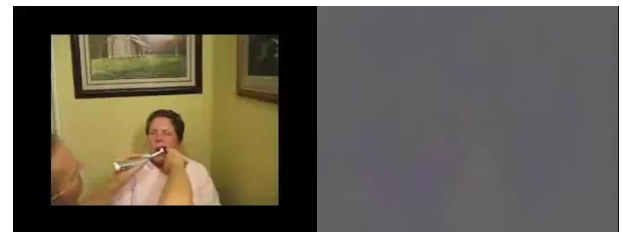
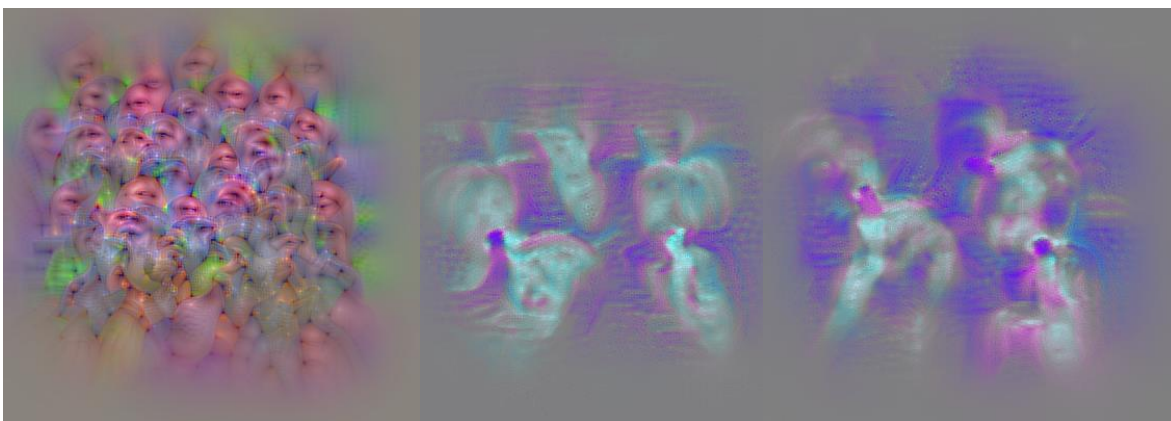
Slow motion

Fast motion



BrushingTeeth 52% accuracy

(33% confused with ShavingBeard)





Revealing idiosyncracies in data

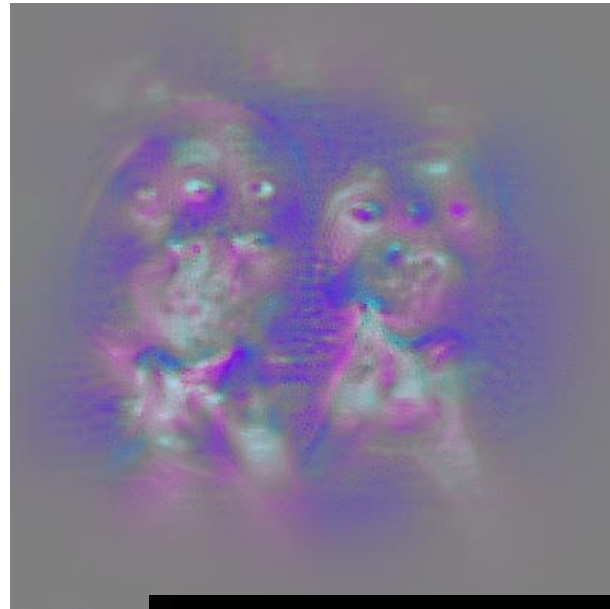
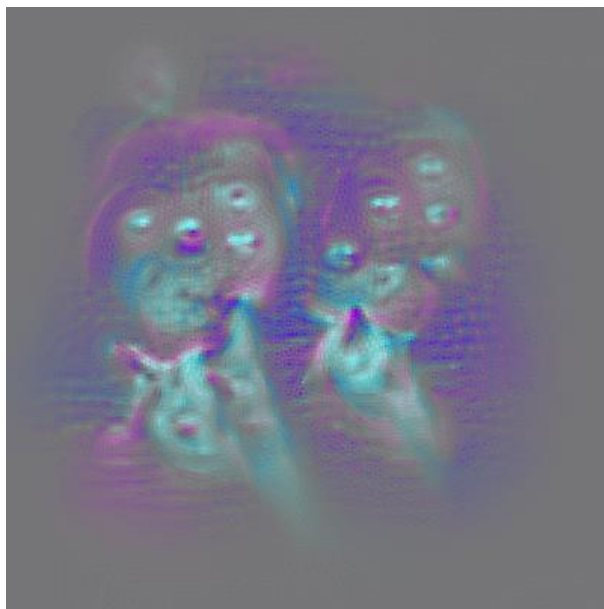
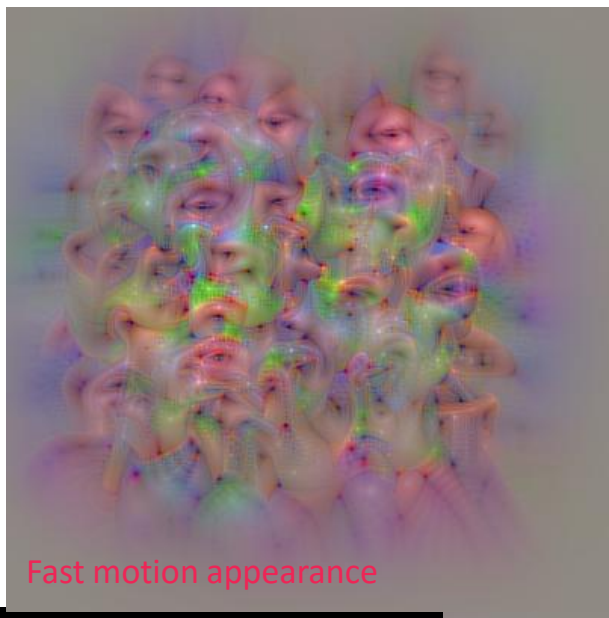
→ “ApplyLipstick”



Appearance

Slow motion

Fast motion



sis 2017





Appearance

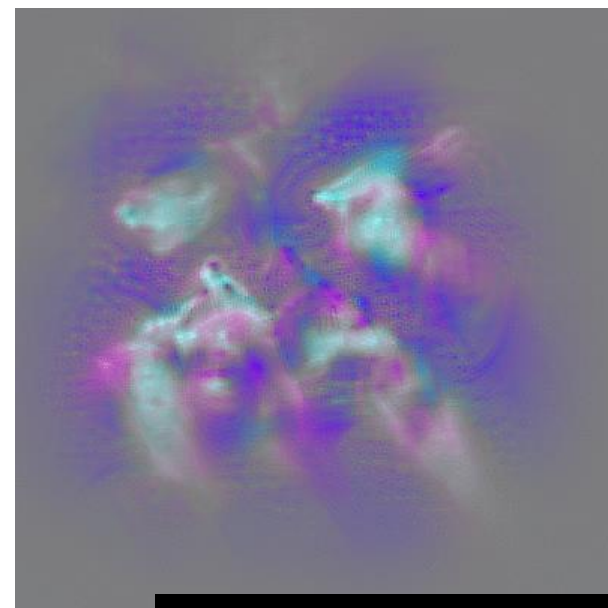
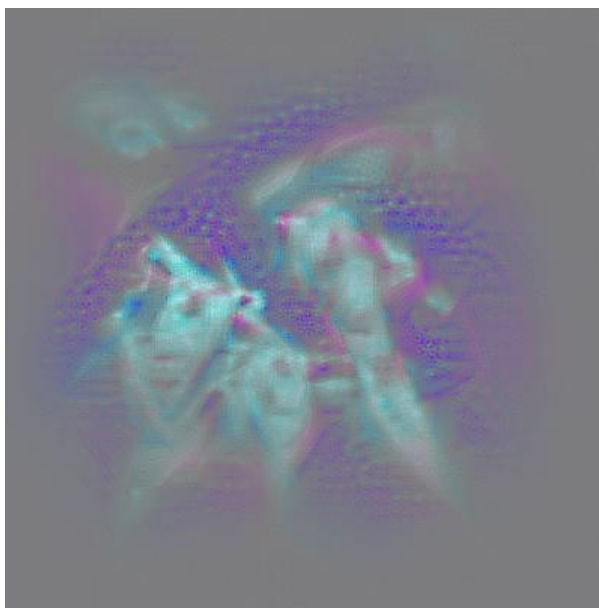
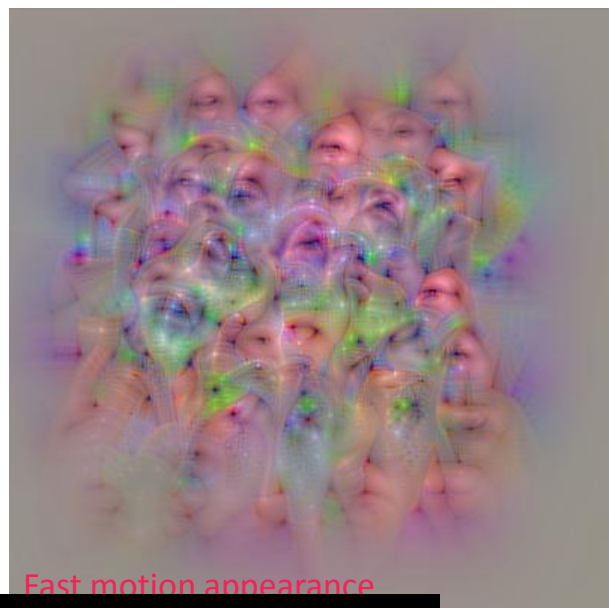
Revealing idiosyncracies in data

→ “ApplyEyeMakeup”



Slow motion

Fast motion

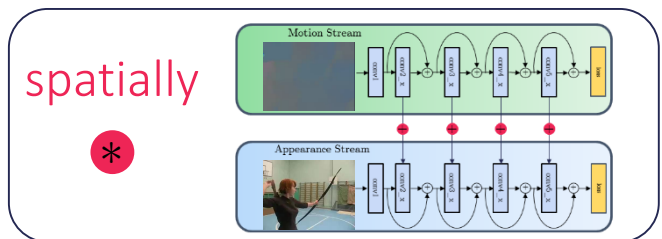


sis 2017



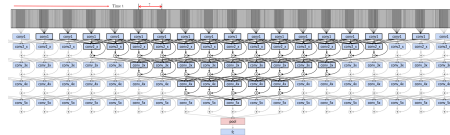
Summary of our insights

We study ways of connecting appearance and motion ConvNets



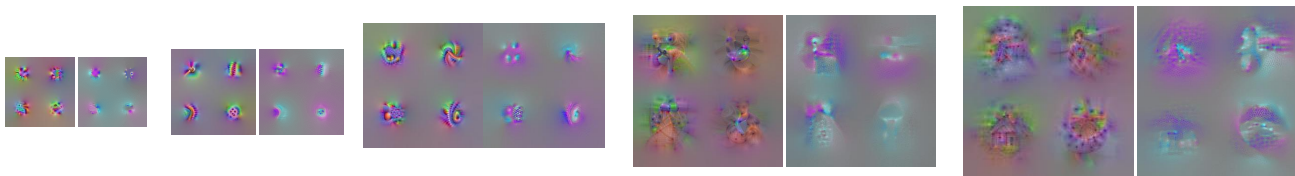
and

temporally *

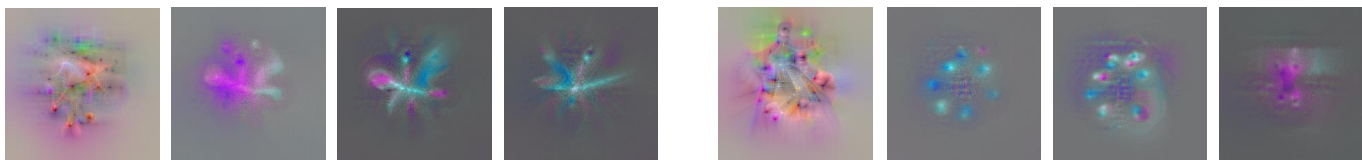


By visualizing the learned representation we find that:

- Early layers show similar spatial structures for appearance and flow



- Higher layer conv-fusion-filters are broadly tuned to multiple speeds and can be specific but also generic across classes



- Class visualizations aid in analyzing system strengths and weaknesses

