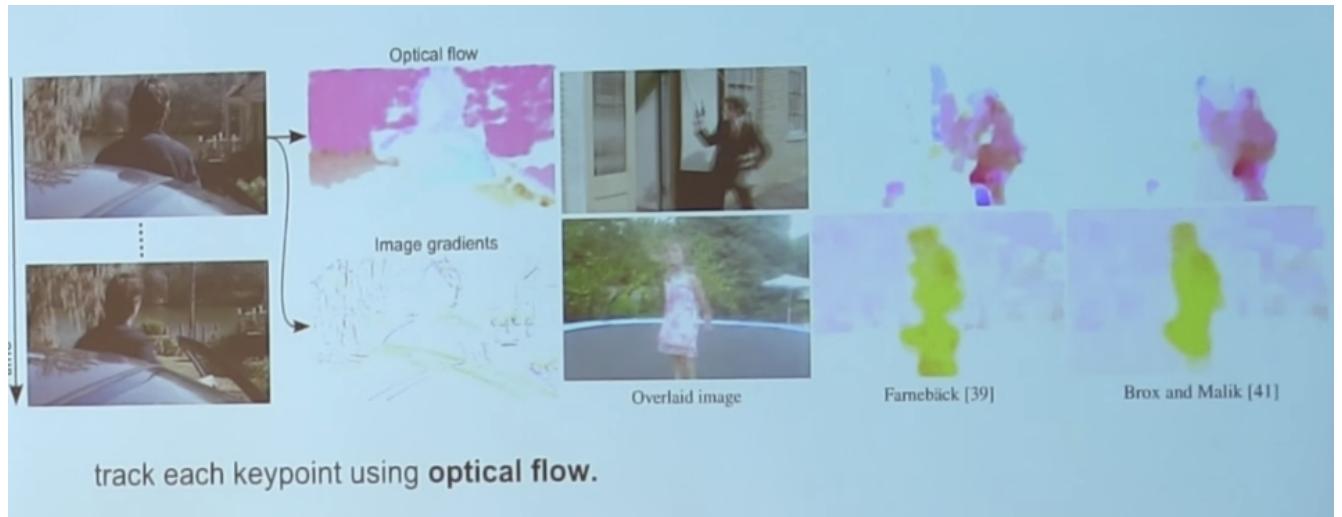
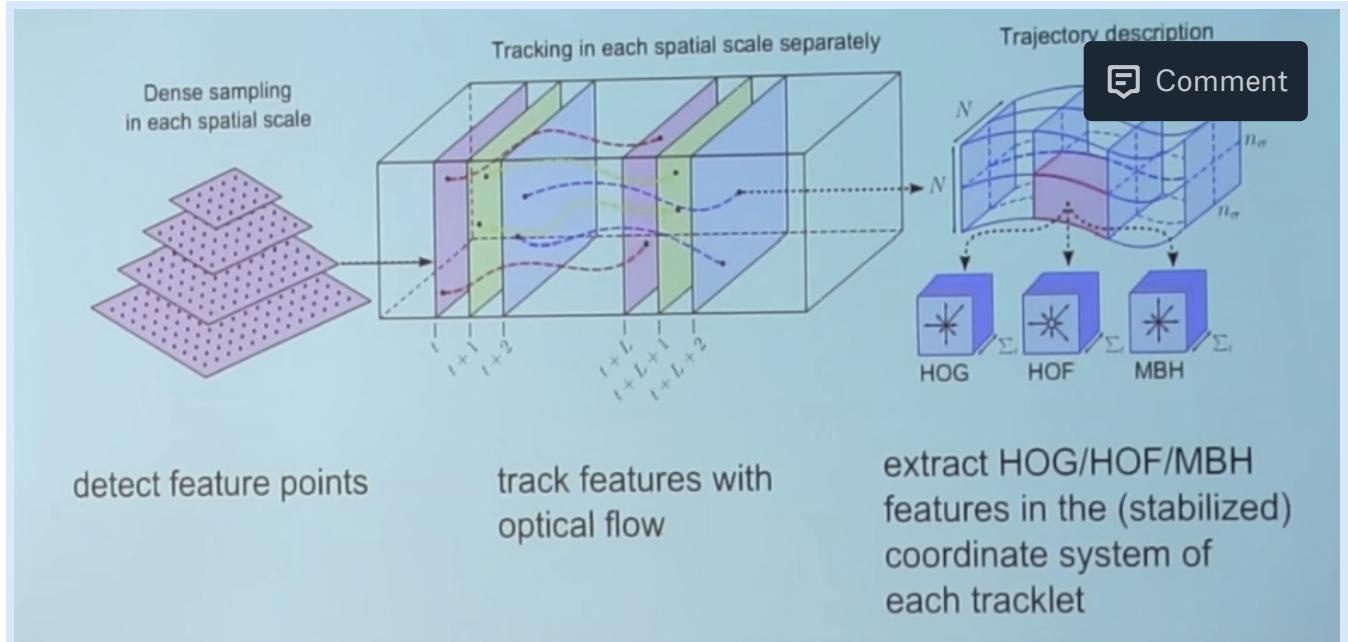


video analysis

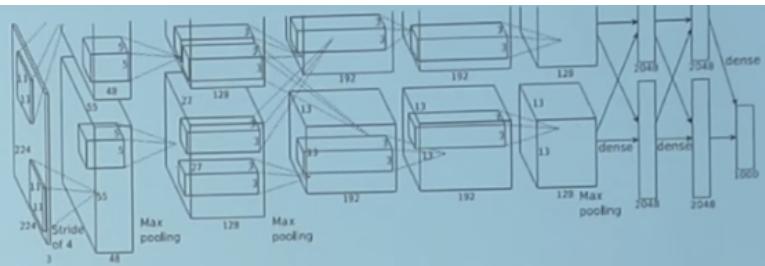
- First, see if you can try a technique from image analysis on your video task. do you really need to use spatio temporal information for your video task? Maybe not!
- before CNNs → they would identify these key feature points to track based on gradients and frequency and other features. they would create frames with the origin at these points and track the points and use that to compute optical flow which is basically a displacement vector of given point across frames.



- then came CNNs

Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

First layer (CONV1): 96 11x11 filters applied at stride 4

=>

Output volume [55x55x96]

Q: What if the input is now a small chunk of video? E.g. [227x227x3x15] ?

A: Extend the convolutional filters in time, perform spatio-temporal convolutions!

E.g. can have 11x11xT filters, where T = 2..15.

- the method in the black square worked best

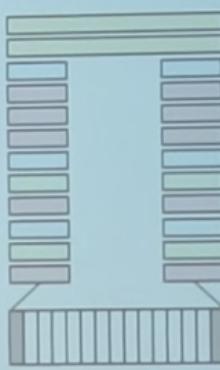
Spatio-Temporal ConvNets

spatio-temporal convolutions;
worked best.

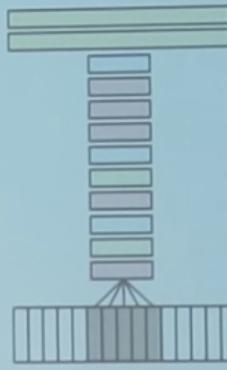
Single Frame



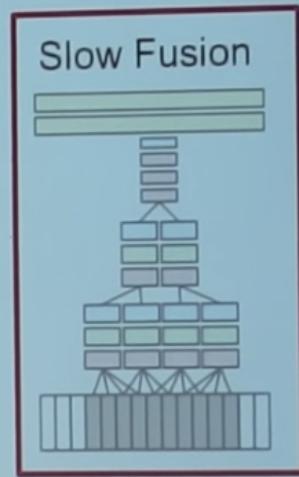
Late Fusion



Early Fusion



Slow Fusion



- maybe you don't need to take advantage of spatio temporal information and can just use vanilla alexnet etc frame by frame

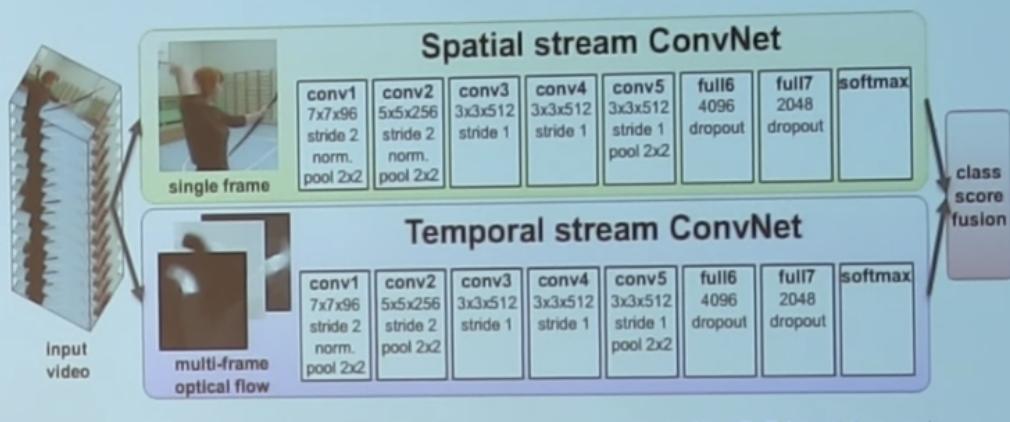
Spatio-Temporal ConvNets

Model	Clip Hit@1	Video Hit@1	Video Hit@5
Feature Histograms + Neural Net	-	55.3	-
Single-Frame	41.1	59.3	77.7
Single-Frame + Multires	42.4	60.0	78.5
Single-Frame Fovea Only	30.0	49.9	72.8
Single-Frame Context Only	38.1	56.0	77.2
Early Fusion	38.9	57.7	76.8
Late Fusion	40.7	59.3	78.7
Slow Fusion	41.9	60.9	80.2
CNN Average (Single+Early+Late+Slow)	41.4	63.9	82.4

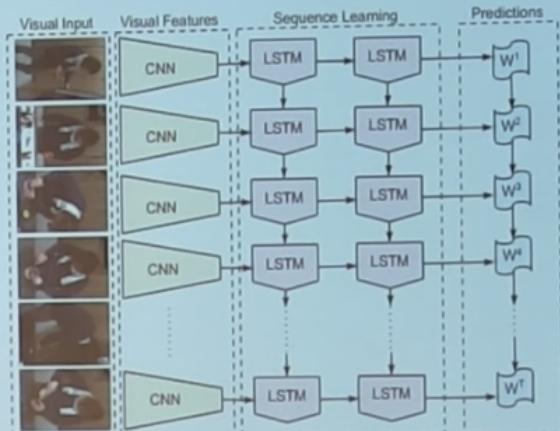
The motion information didn't add all that much...

- compute optical flow image and use that as a feature

Spatio-Temporal ConvNets



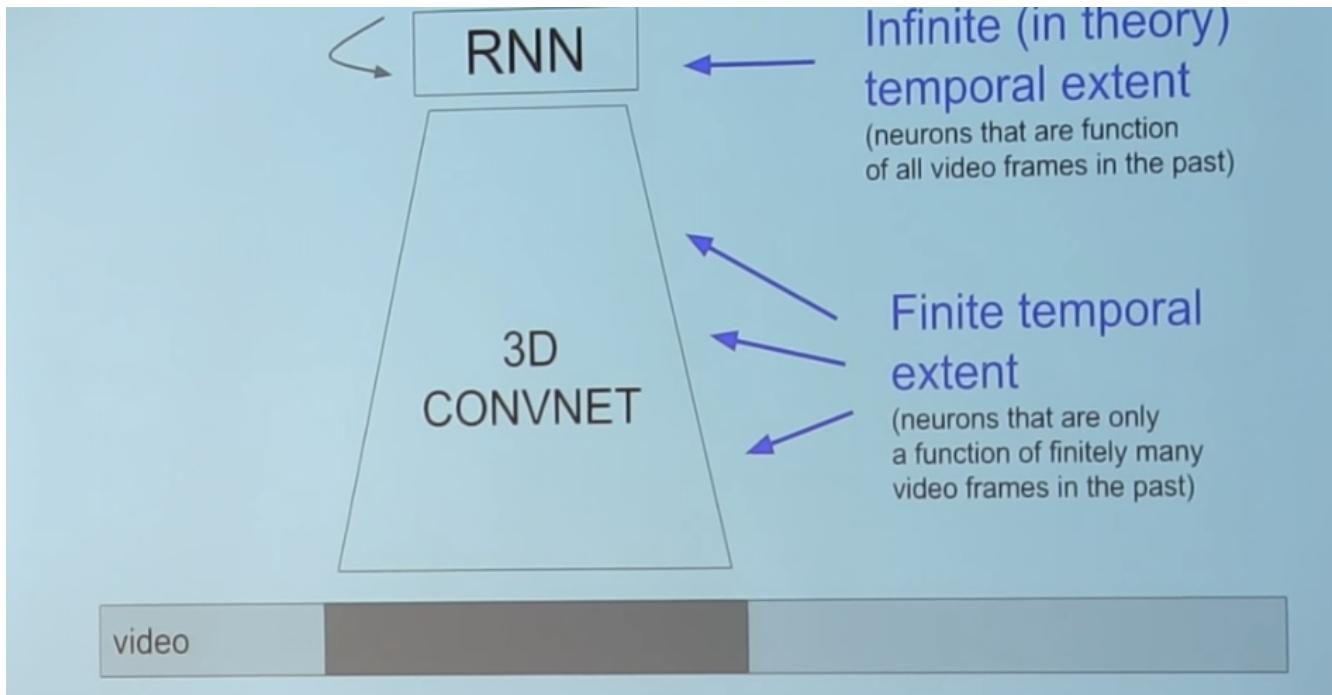
Long-time Spatio-Temporal ConvNets



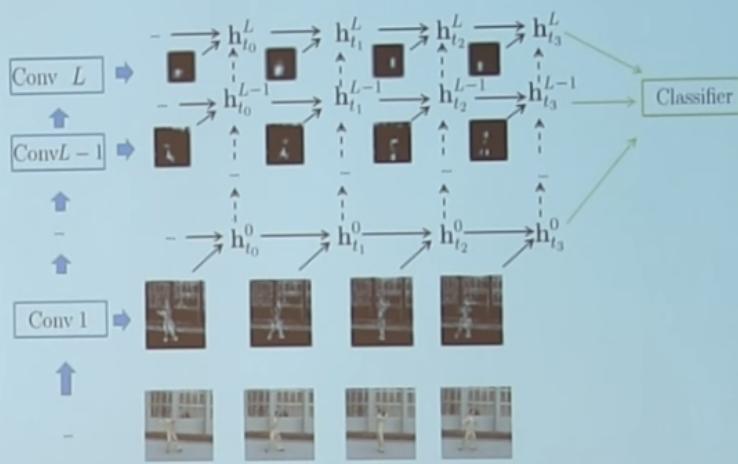
Summary so far

We looked at two types of architectural patterns:

1. Model temporal motion locally (3D CONV)
2. Model temporal motion globally (LSTM / RNN)
 - + Fusions of both approaches at the same time.



Long-time Spatio-Temporal ConvNets



Beautiful:
All neurons in the ConvNet are recurrent.

$$\begin{aligned} z_t^l &= \sigma(W_z^l * x_t^l + U_z^l * h_{t-1}^l), \\ r_t^l &= \sigma(W_r^l * x_t^l + U_r^l * h_{t-1}^l), \\ \tilde{h}_t^l &= \tanh(W^l * x_t^l + U * (r_t^l \odot h_{t-1}^l)), \\ h_t^l &= (1 - z_t^l)h_{t-1}^l + z_t^l \tilde{h}_t^l, \end{aligned}$$

Only requires (existing) 2D CONV routines. No need for 3D spatio-temporal CONV.

Long-time Spatio-Temporal ConvNets

