

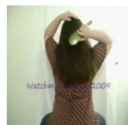


# Temporal Models for Video Processing

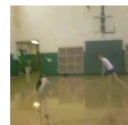


# Datasets

- HMDB-51: 7000 videos of 51 actions
- UFC 101: 13,320 videos of 101 actions
- Kinetics: Up to 650,000 videos, up to 700 actions



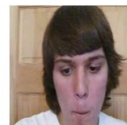
brush hair



cartwheel



catch



chew



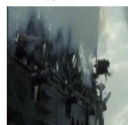
clap



climb



climb stairs



dive



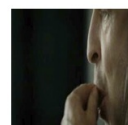
draw sword



dribble



drink



eat



fall floor



fencing



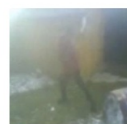
flic flac



golf



hand stand



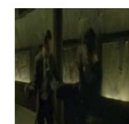
hit



hug



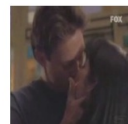
jump



kick



kick ball



kiss



laugh



pick



pour

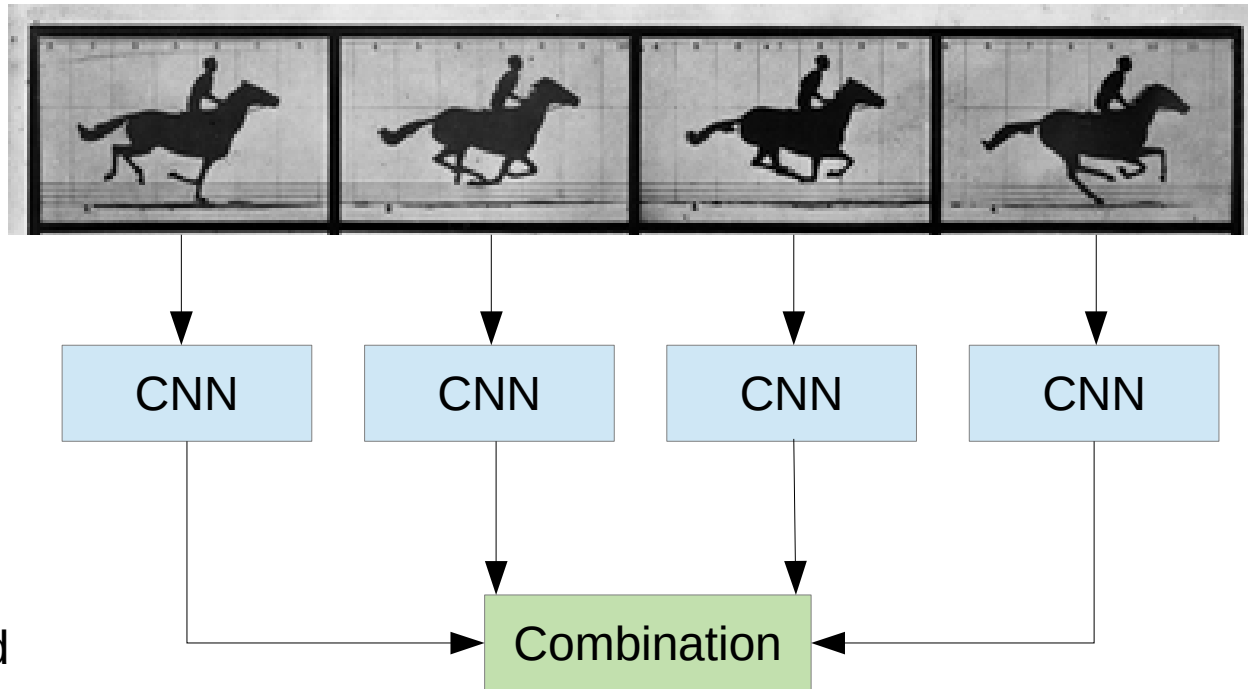


pullup



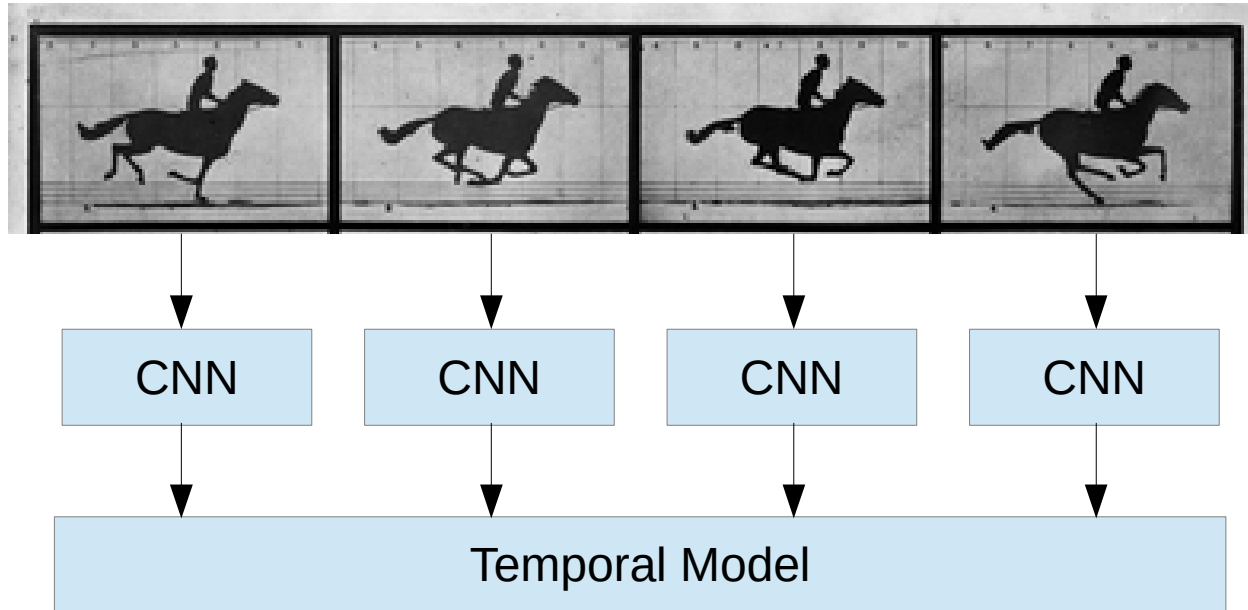
punch

# Approach 1: Unordered Frames



Pretty good  
baseline

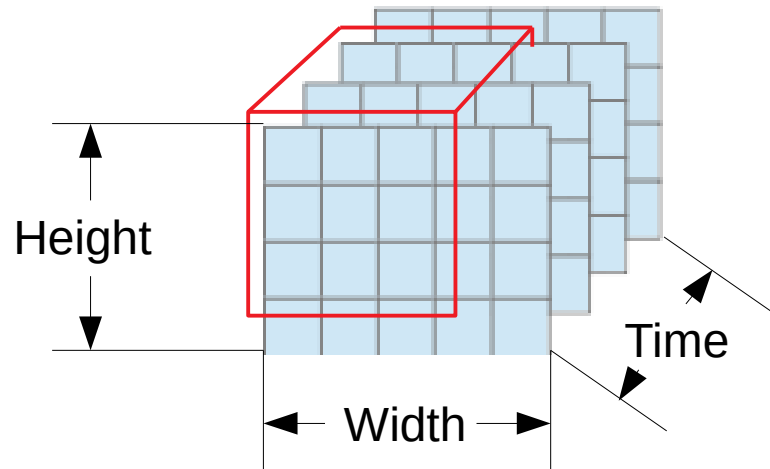
# Approach 2: Frames + Global Model



Tends to overfit, hard to train

# 3D Convolutions

Convolution across both space and time



# 3D Convolutions

Input  $X \in \mathbb{R}^{C_i \times \underline{T} \times H \times W}$

Kernel  $W \in \mathbb{R}^{C_o \times C_i \times \underline{t} \times h \times w}$

Bias  $b \in \mathbb{R}^{C_o}$

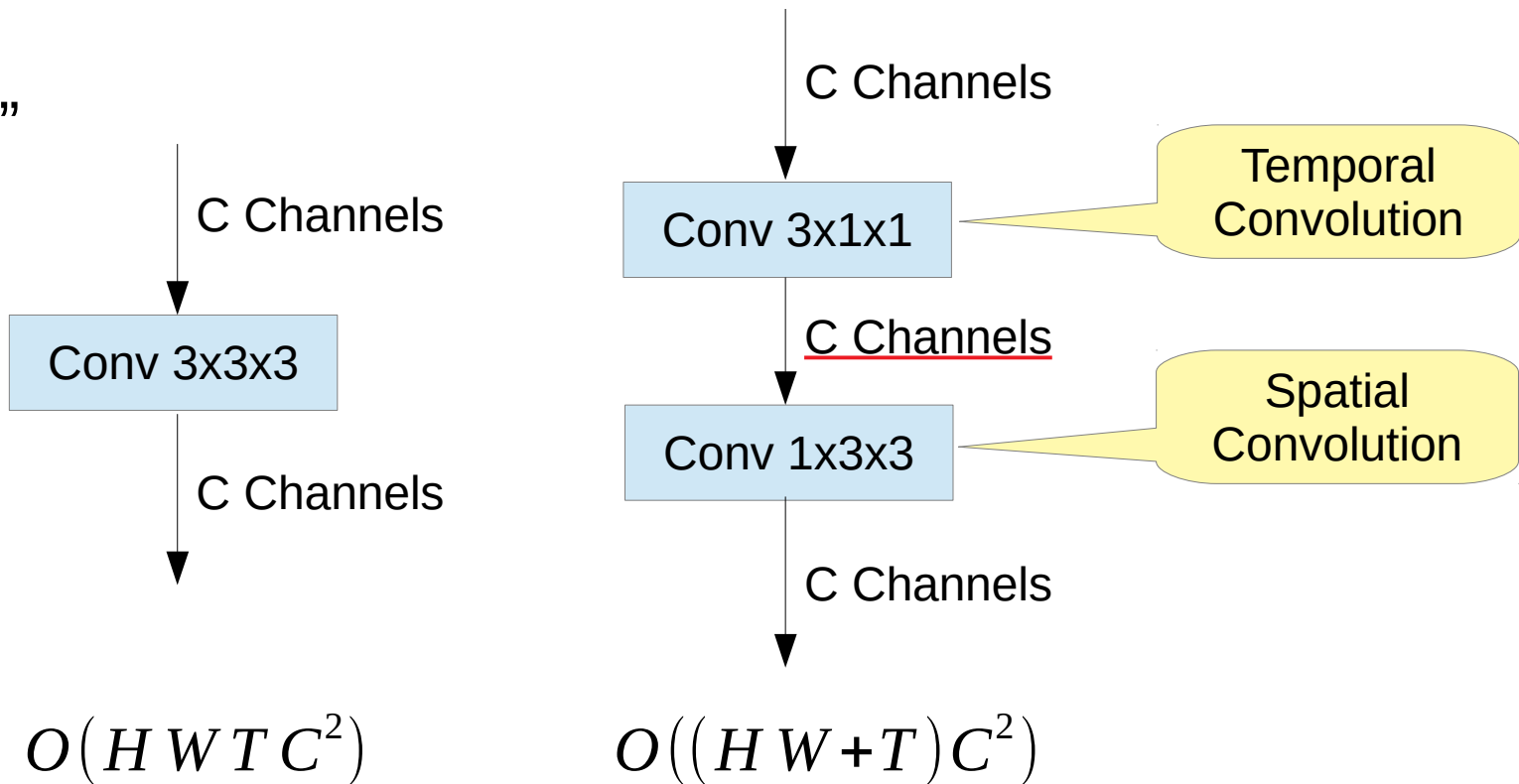
Output  $Z \in \mathbb{R}^{C_o \times \left( \frac{T-t+2p_t}{s_t} + 1 \right) \times \left( \frac{H-h+2p_h}{s_h} + 1 \right) \times \left( \frac{W-w+2p_w}{s_w} + 1 \right)}$

$$Z_{c, \underline{d}, a, b} = b_c + \sum_{l=0}^{C_i} \sum_{\underline{i}=0}^t \sum_{j=0}^h \sum_{k=0}^w X_{k, \underline{d+i}, a+j, b+k} W_{c, l, \underline{i}, j, k}$$

Very slow!

# Factorized 3D Convolutions

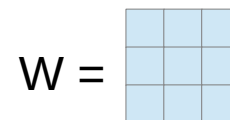
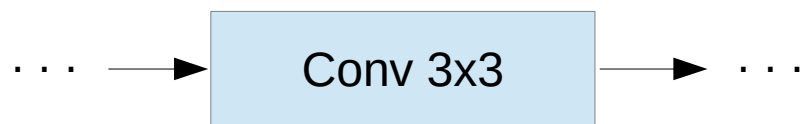
“2+1D”



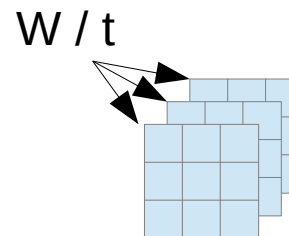
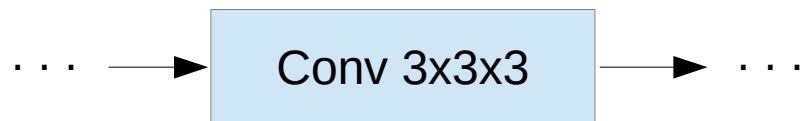


# I3D

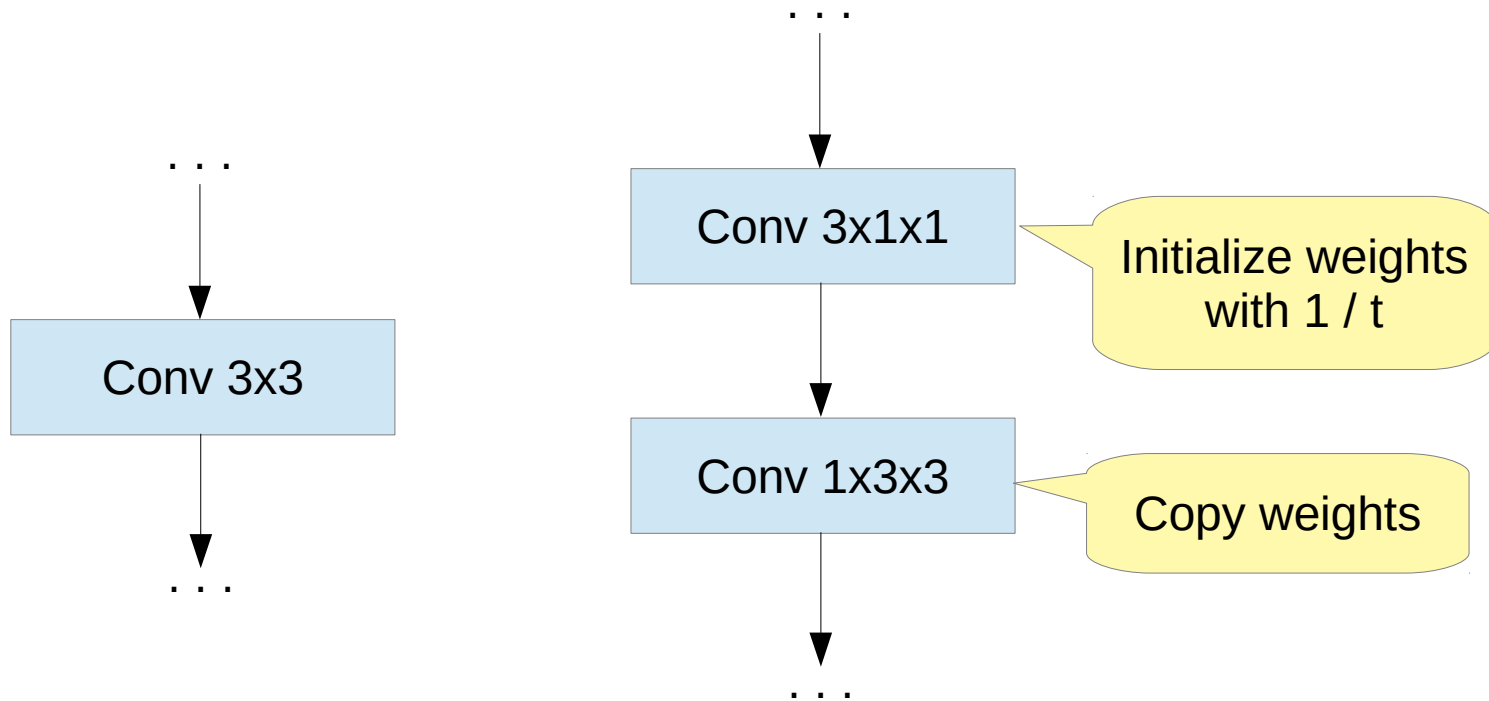
- Pre-train a network on ImageNet



- “Inflate” some 2D convolutions to 3D



# I2+1D



# Open Problem: What Tasks Should We Care About?

- Vision tasks are often proxies or initial steps in other applications
  - But they often don't really capture the downstream task.
- Vision is good as a test bed.
  - New architectures
  - Pre-training

