

# Video Classification

**Danna Gurari**

University of Colorado Boulder  
Fall 2021



# Review

- Last week's lecture topic:
  - Salient object detection
- Assignments (Canvas)
  - Reading assignment due earlier today
  - Reading assignments due Wednesday and next week
  - Final project proposal due in just over 3 weeks
- Questions?

# Final Project Requirements

- Described on the course website
  - <https://home.cs.colorado.edu/~DrG/Courses/RecentAdvancesInComputerVision/FinalProject.html>
- Multiple milestones
  - Project proposal
  - Project outline
  - Final project presentation
  - Peer evaluation
  - Final report

# Video Classification: Today's Topics

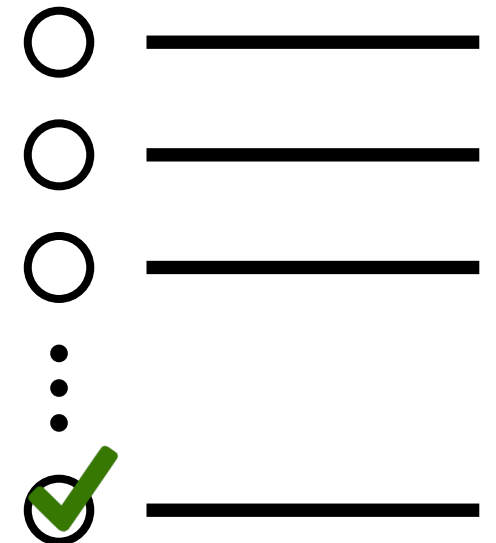
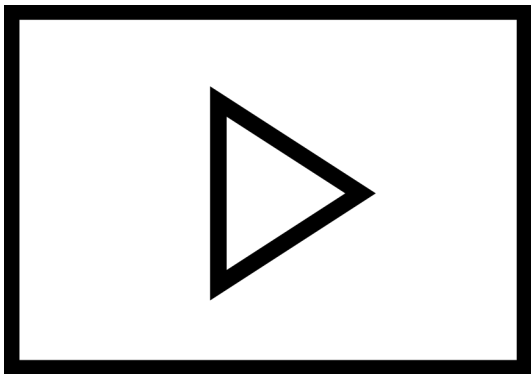
- Problem
- Applications
- Datasets
- Evaluation metric
- Computer vision models

# Video Classification: Today's Topics

- Problem
- Applications
- Datasets
- Evaluation metric
- Computer vision models

# Definition

- Assign a video a label from a set of categories; typically, multiple choice but also can be multiple labels
  - e.g., activity or topical themes



# Video Classification/Localization: Today's Topics

- Problem
- Applications
- Datasets
- Evaluation metric
- Computer vision models

# Video Search

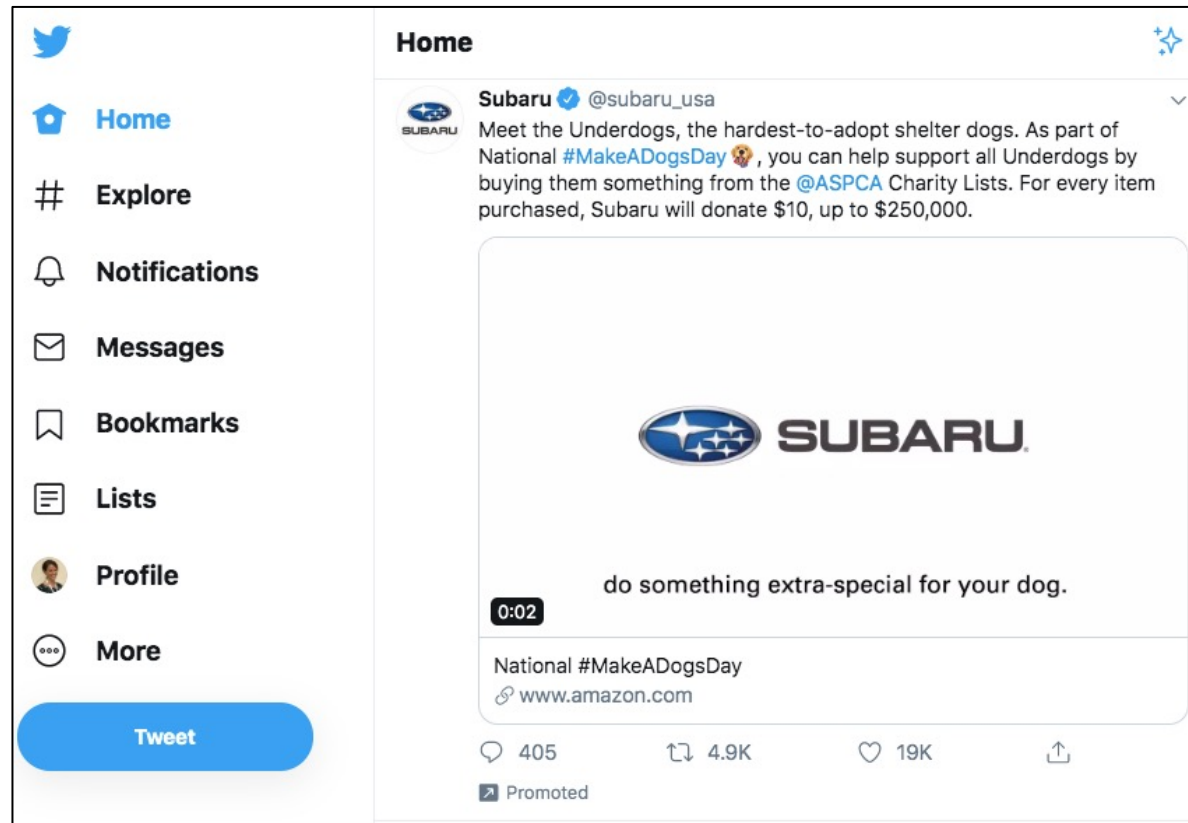
The image shows a screenshot of the YouTube homepage. At the top, there is a search bar with the text "Search" and a magnifying glass icon. To the right of the search bar are icons for a camera, a grid, and a vertical ellipsis, followed by a "SIGN IN" button. On the left side, there is a navigation menu with the following items: "Home" (selected), "Trending", "Subscriptions", "Library", and "History". Below the navigation menu, there is a sign-in prompt: "Sign in to like videos, comment, and subscribe." with a "SIGN IN" button. Underneath that is a section titled "BEST OF YOUTUBE" with icons and links for "Music", "Sports", "Gaming", and "Movies". The main content area features a large "YouTube TV" banner with the text "Try it free" and a play button icon. To the right of the banner is a "Watch the World Series" advertisement for YouTube TV with a "TRY IT FREE" button. Below the banner is a "Trending" section with five video thumbnails. Each thumbnail includes a title, channel name, and view count.

Thumbnail 1	Thumbnail 2	Thumbnail 3	Thumbnail 4	Thumbnail 5
<b>You Can't Con a Con Artist If You're Also a Con Artist - Ke...</b>	<b>Patriots vs. Jets Week 7 Highlights   NFL 2019</b>	<b>TECHNIQUE CRITIQUE S1 • E14</b>	<b>Unexpected Trick Shots   Dude Perfect</b>	<b>\$3.50 Soup Vs. \$29 Soup • Taiwan</b>
Key & Peele 2M views • 18 hours ago	NFL 1M views • 11 hours ago	WIRED 790K views • 22 hours ago	Dude Perfect 6M views • 16 hours ago	BuzzFeedVideo 2.1M views • 1 day ago

300 hours of video uploaded every minute (<https://merchdope.com/youtube-stats/>)



# Social Media Recommendations



“An estimated 12 million micro-videos are posted to Twitter each day. The number of microvideos produced surpasses the total inventory of YouTube every 3 months”

- “The Open World of Micro-Videos; Nguyen et al.; [https://www.ics.uci.edu/~fowlkes/papers/nrfr\\_bigvision.pdf](https://www.ics.uci.edu/~fowlkes/papers/nrfr_bigvision.pdf)

# Video Organization



Lists search results based on your collection of videos (spanning YouTube, news, movies, and more) in one list

# Automatically Remove Objectionable Content



## Nudity or sexual content

YouTube is not for pornography or sexually explicit content. If this describes your video, even if it's a video of yourself, don't post it on YouTube. Also, be advised that we work closely with law enforcement and we report child exploitation. [Learn more](#)



## Harmful or dangerous content

Don't post videos that encourage others to do things that might cause them to get badly hurt, especially kids. Videos showing such harmful or dangerous acts may get age-restricted or removed depending on their severity. [Learn more](#)



## Hateful content

Our products are platforms for free expression. But we don't support content that promotes or condones violence against individuals or groups based on race or ethnic origin, religion, disability, gender, age, nationality, veteran status, or sexual orientation/gender identity, or whose primary purpose is inciting hatred on the basis of these core characteristics. This can be a delicate balancing act, but if the primary purpose is to attack a protected group, the content crosses the line.

[Learn more](#)



## Violent or graphic content

It's not okay to post violent or gory content that's primarily intended to be shocking, sensational, or gratuitous. If posting graphic content in a news or documentary context, please be mindful to provide enough information to help people understand what's going on in the video. Don't encourage others to commit specific acts of violence. [Learn more](#)

And more listed here: <https://www.youtube.com/about/policies/#community-guidelines>

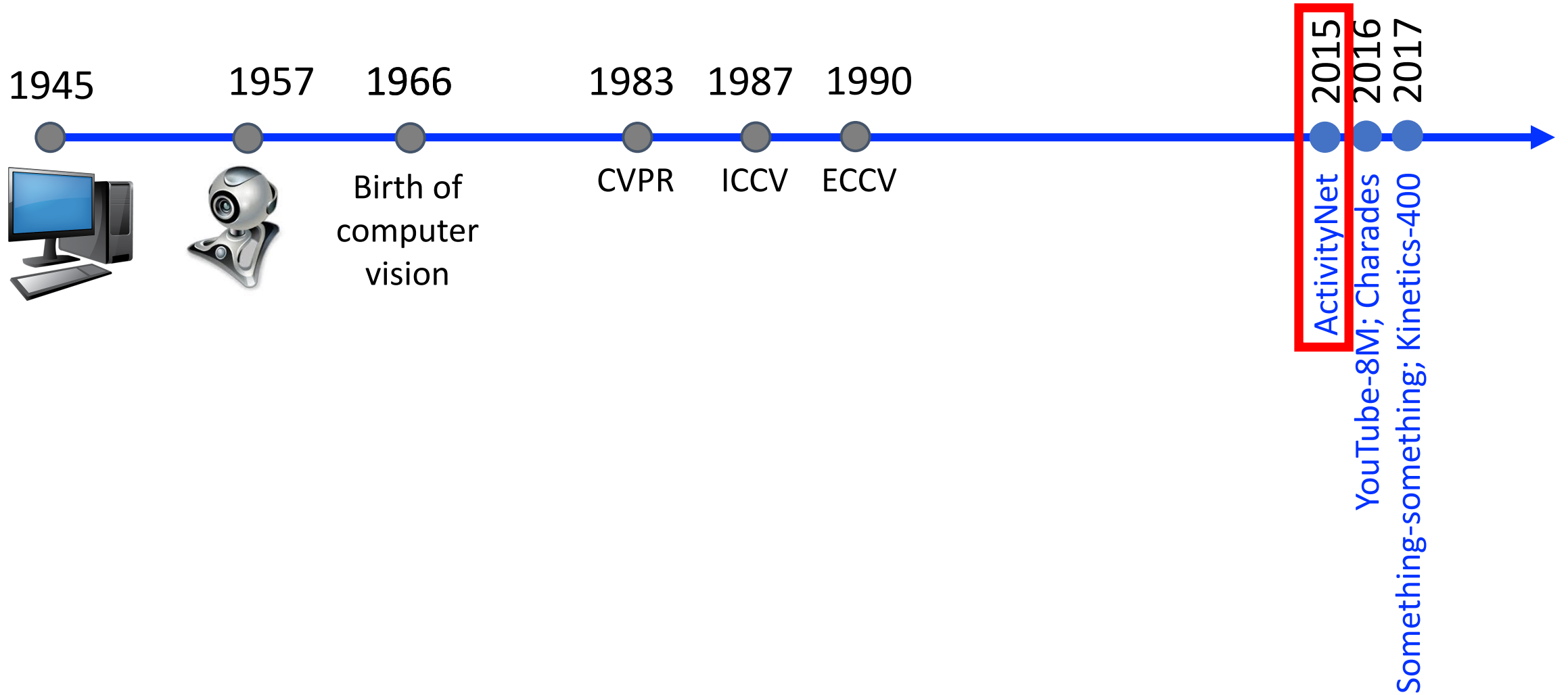
# Applications

For what other applications might video classification be useful?

# Video Classification: Today's Topics

- Problem
- Applications
- **Datasets**
- Evaluation metric
- Computer vision models

# Datasets



# ActivityNet

Focus on activities that humans spend most of their time doing in their lives

# ActivityNet

## 1. Category Selection

- \* American Time Use Survey (ATUS) created by the Department of Labor organizes activities according to:
  - social interactions
  - where activity usually occurs
- \* Authors selected 203 from the 2000+ activities in ATUS:
  - 7 top-level categories:
  - Personal Care, Eating and Drinking, Household, Working,...*
  - 4-level hierarchy





# ActivityNet

## 1. Category Selection

- \* American Time Use Survey (ATUS) created by the Department of Labor organizes activities according to:
  - social interactions
  - where activity usually occurs
- \* Authors selected 203 from the 2000+ activities in ATUS:
  - 7 top-level categories: *Personal Care, Eating and Drinking, Household, Working,...*
  - 4-level hierarchy

## 2. Video Collection



## 3. Video Verification

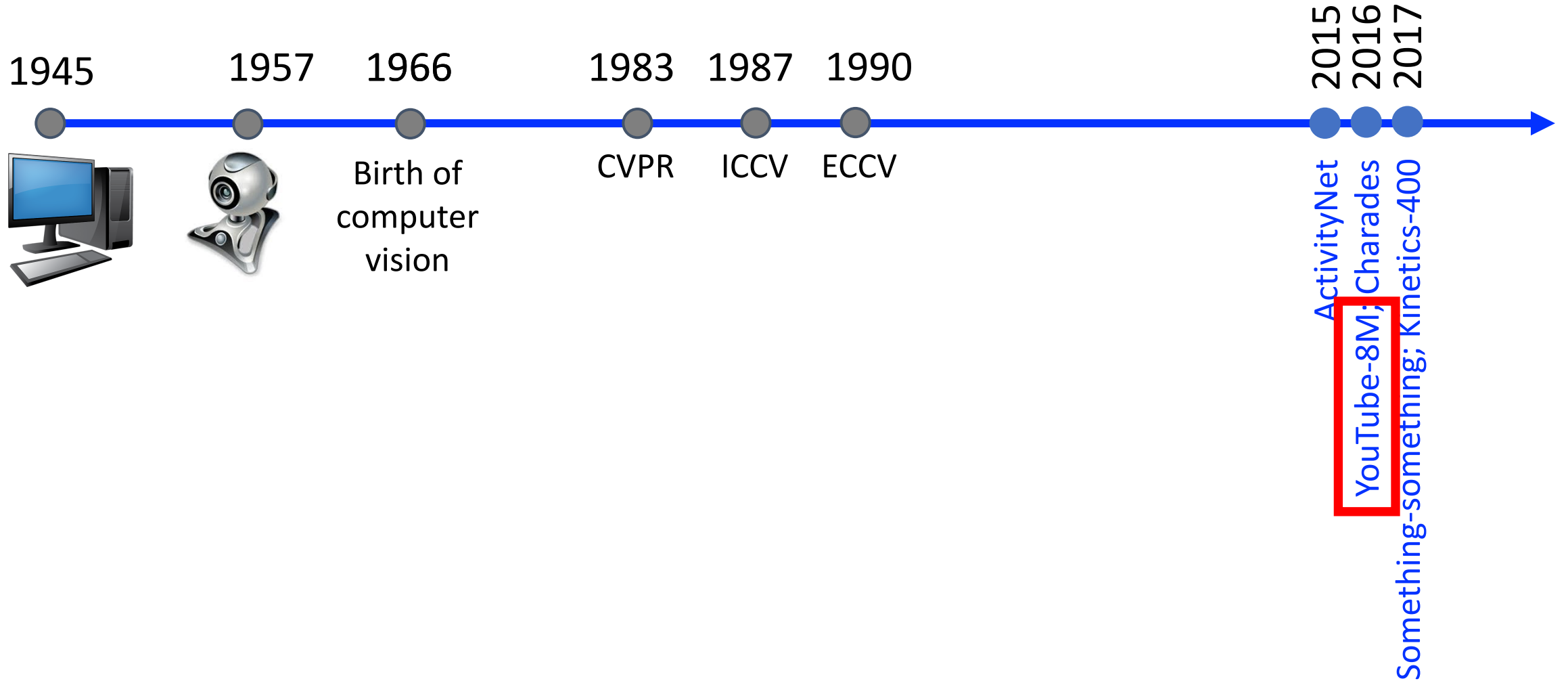
- \* "Expert" AMT workers verify presence of activity in each video
- \* Honey pot tasks introduced to assess trust of each crowd worker's work

# ActivityNet Workshop

The image shows a screenshot of the ActivityNet website. At the top left is the ActivityNet logo, which consists of a red square with a white silhouette of a person jumping, followed by the text 'ACTIVITYNET' in bold black letters and 'Large Scale Activity Recognition Challenge' in smaller black text below it. To the right of the logo is a navigation menu with the following items: 'HOME' (in red), 'PEOPLE', 'CHALLENGE', 'PROGRAM', 'DATES', 'EVALUATION', and 'CONTACT'. Further to the right is a blue rectangular button with the text 'CVPR VIRTUAL JUNE 19-25' in white. Below the navigation is a large banner image showing a busy indoor scene, possibly a train station or a public event, with many people. Overlaid on this banner is the text 'International Challenge on Activity Recognition (ActivityNet)' in large white font, and 'CVPR 2021 Workshop' in slightly smaller white font below it. A red diagonal shape is visible in the bottom left corner of the banner area.

<http://activity-net.org/challenges/2021/>

# Datasets



# YouTube-8M

Largest multi-label video classification dataset for determining the key topical themes of the video

- ~8 million videos of over 500,000 hours
- 4,800 classes spanning “activities (sports, games, hobbies), objects (autos, food, products), scenes (travel), and events”

# YouTube-8M

## 1. Category Selection

\* Starting point: 50,000 video topics from a knowledge graph (Freebase); e.g., people, places

\* Kept ~10,000 topics that most of 3 humans indicated are visually distinguishable and do not require domain expertise to recognize

\* Reduced to categories that are popular: 1,000+ views, > 120 secs, < 500 secs, and >= 200 videos

Manual pruning task:

Entity Name	Entity URL	Entity Description
Thunderstorm	<a href="http://www.freebase.com/m/0jb2l">http://www.freebase.com/m/0jb2l</a>	A thunderstorm, also known as an electrical storm, a lightning storm, or a thundershower, is a type of storm characterized by the presence of lightning and its acoustic effect on the Earth's atmosphere known as thunder. The meteorologically assigned cloud type associated with the thunderstorm is the cumulonimbus. Thunderstorms are usually accompanied by strong winds, heavy rain and sometimes snow, sleet, hail, or no precipitation at all...

How difficult is it to identify this entity in images or videos (without audio, titles, comments, etc)?

- 1. Any layperson could
- 2. Any layperson after studying examples, wikipedia, etc could
- 3. Experts in some field can
- 4. Not possible without non-visual knowledge
- 5. Non-visual

# YouTube-8M

## 1. Category Selection

- \* Starting point: 50,000 video topics from a knowledge graph (Freebase); e.g., people, places
- \* Kept ~10,000 topics that most of 3 humans indicated are visually distinguishable and do not require domain expertise to recognize
- \* Reduced to categories that are popular: 1,000+ views, > 120 secs, < 500 secs, and  $\geq$  200 videos

## 2. Video Collection



## 3. Label Verification

- \* 3 humans rate labels for 8000 images
- \* With respect to the human raters
  - 78.8% precision
  - 14.5% recall
- \* Inter-rater agreement: ~80%

# YouTube-8M Challenge & Annual Workshop

YouTube | 8M

Dataset

Explore

Download

Workshop

About

## Updated Dataset

YouTube-8M Segments was released in June 2019 with segment-level annotations. Human-verified labels on about 237K segments and 1000 classes are collected from the validation set of the YouTube-8M dataset. Each video will again come with time-localized features so classifier predictions can be made at segment-level granularity.

YouTube-8M was updated in May 2018 to include higher-quality, more topical annotations, and to clean up the annotation vocabulary. A number of low-frequency or low-quality labels and associated videos were removed, resulting in a smaller but higher-quality dataset (5.6M videos, 3862 classes). Additionally, the video IDs in the TensorFlow Record files have been anonymized, and the mapping to the real YouTube IDs will be periodically updated to exclude any videos that have been subsequently deleted (while preserving their anonymized features).

Dataset versions:

1. Jun 2019 version (current): 230K human-verified segment labels, 1000 classes, 5 segments/video
2. May 2018 version (current): 6.1M videos, 3862 classes, 3.0 labels/video, 2.6B audio-visual features
3. Feb 2017 version (deprecated): 7.0M videos, 4716 classes, 3.4 labels/video, 3.2B audio-visual features
4. Sep 2016 version (deprecated): 8.2M videos, 4800 classes, 1.8 labels/video, 1.9B visual-only features

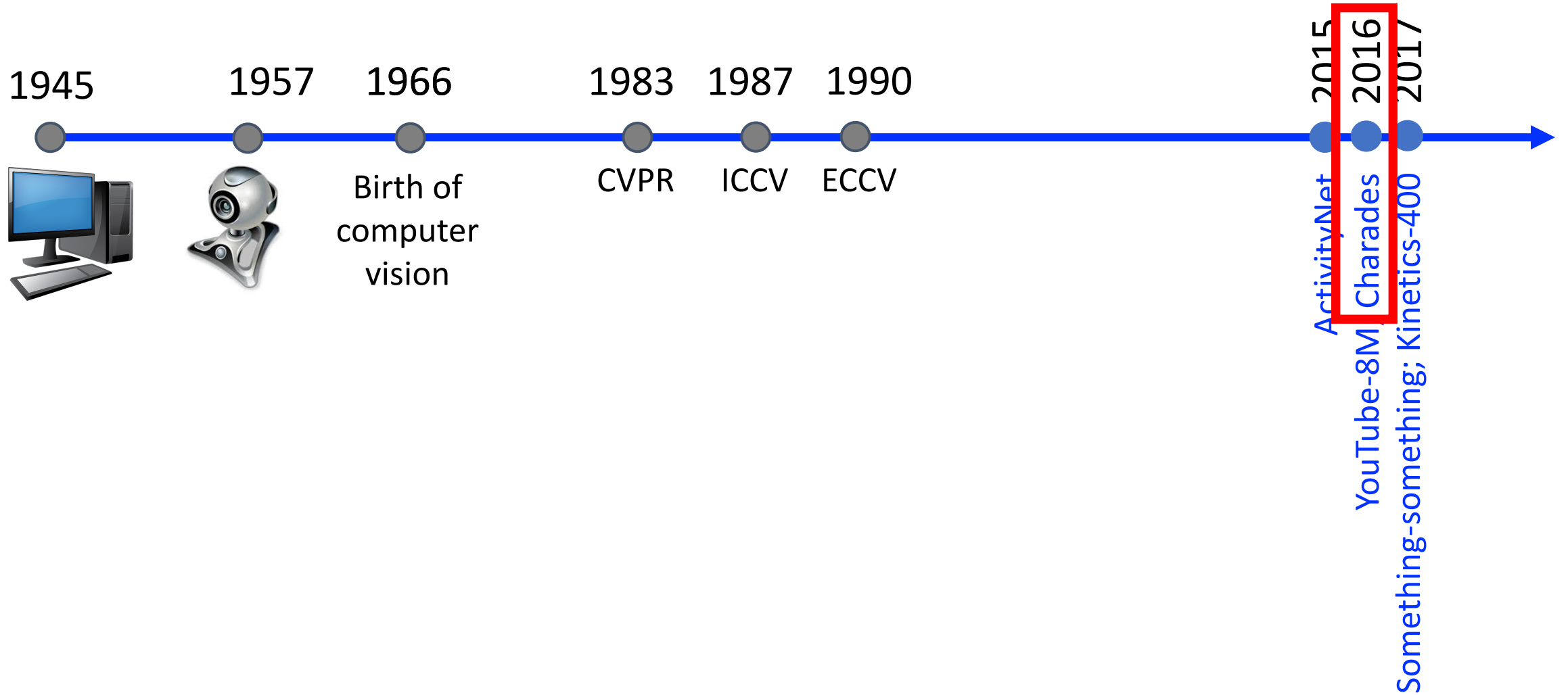
2019

2018

2017

<https://research.google.com/youtube8m/workshop2019/>

# Datasets





# Charades

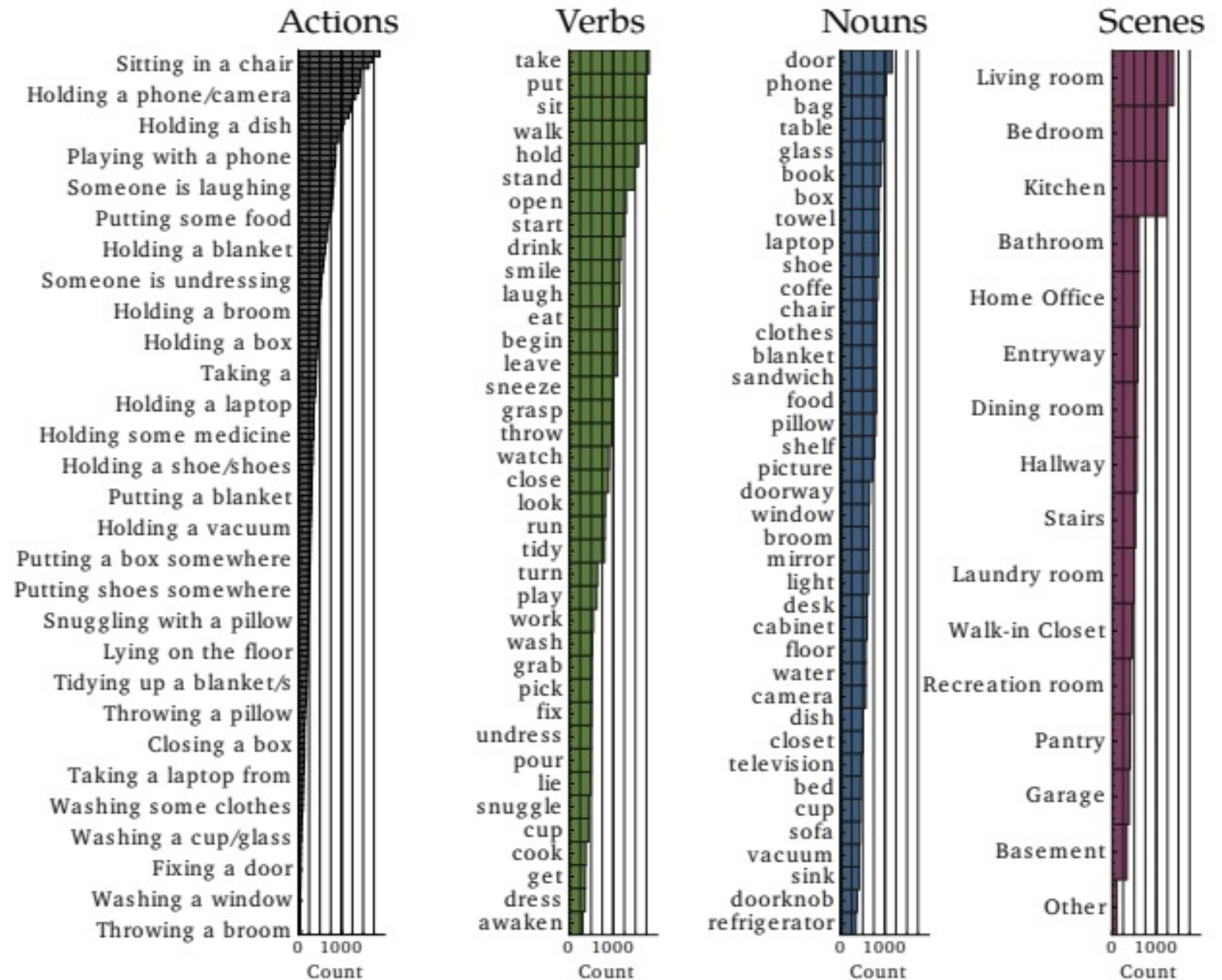
Collection of “boring” videos reflecting daily lives

- ~9,848 videos with average length of 30 seconds
- Activities of 267 people from three continents
- Annotations include action labels and classes of interacted objects

# Charades

## 1. Video Script Generation

- \* Authors identified 15 indoor scenes in residential homes (e.g., living room, home office)
- \* Most common nouns and verbs in these scenes analyzed from 549 movie scripts resulting in 40 objects and 30 actions
- \* Crowd workers generated scripts describing commonplace, realistic activities that involve 2 objects & 2 actions (given a scene, 5 objects, & 5 actions)



# Charades

## 1. Video Script Generation

- \* Authors identified 15 indoor scenes in residential homes (e.g., living room, home office)
- \* Most common nouns and verbs in these scenes analyzed from 549 movie scripts resulting in 40 objects and 30 actions
- \* Crowd workers generated scripts describing commonplace, realistic activities that involve 2 objects & 2 actions (given a scene, 5 objects, & 5 actions)

## 2. Video Collection

- \* Crowd workers recruited to record 30s videos of them executing the scripts

### Demo of videos

<https://www.youtube.com/watch?v=x9AhZLDkbyc>

# Charades

## 1. Video Script Generation

- \* Authors identified 15 indoor scenes in residential homes (e.g., living room, home office)
- \* Most common nouns and verbs in these scenes analyzed from 549 movie scripts resulting in 40 objects and 30 actions
- \* Crowd workers generated scripts describing commonplace, realistic activities that involve 2 objects & 2 actions (given a scene, 5 objects, & 5 actions)

## 2. Video Collection

- \* Crowd workers recruited to record 30s videos of them executing the scripts

## 3. Category Selection

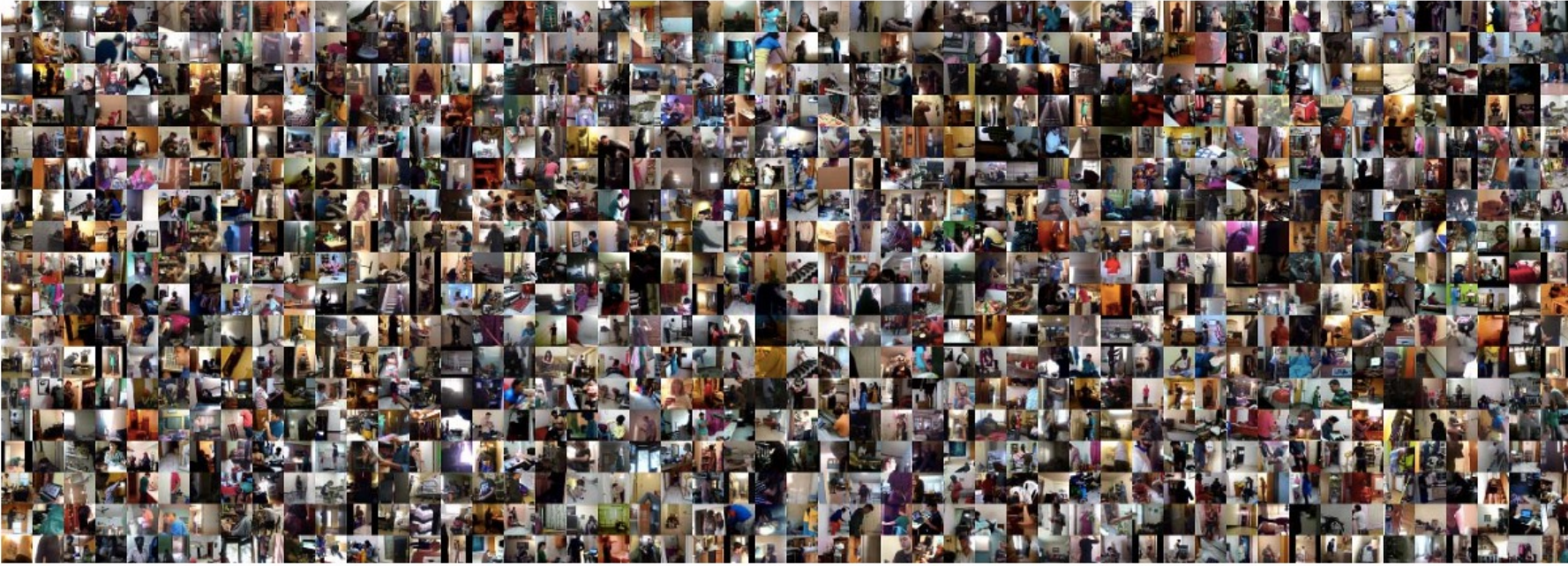
- \* AMT workers recruited to watch each video and create a description
- \* Automatically identified “interacted objects” mentioned both in script & description
- \* 150 actions chosen following crowd worker verification

# Charades Challenge & Annual Workshop

CVPR 2017 Workshop on Visual Understanding Across Modalities Home THOR Charades TQA

## Charades Challenge

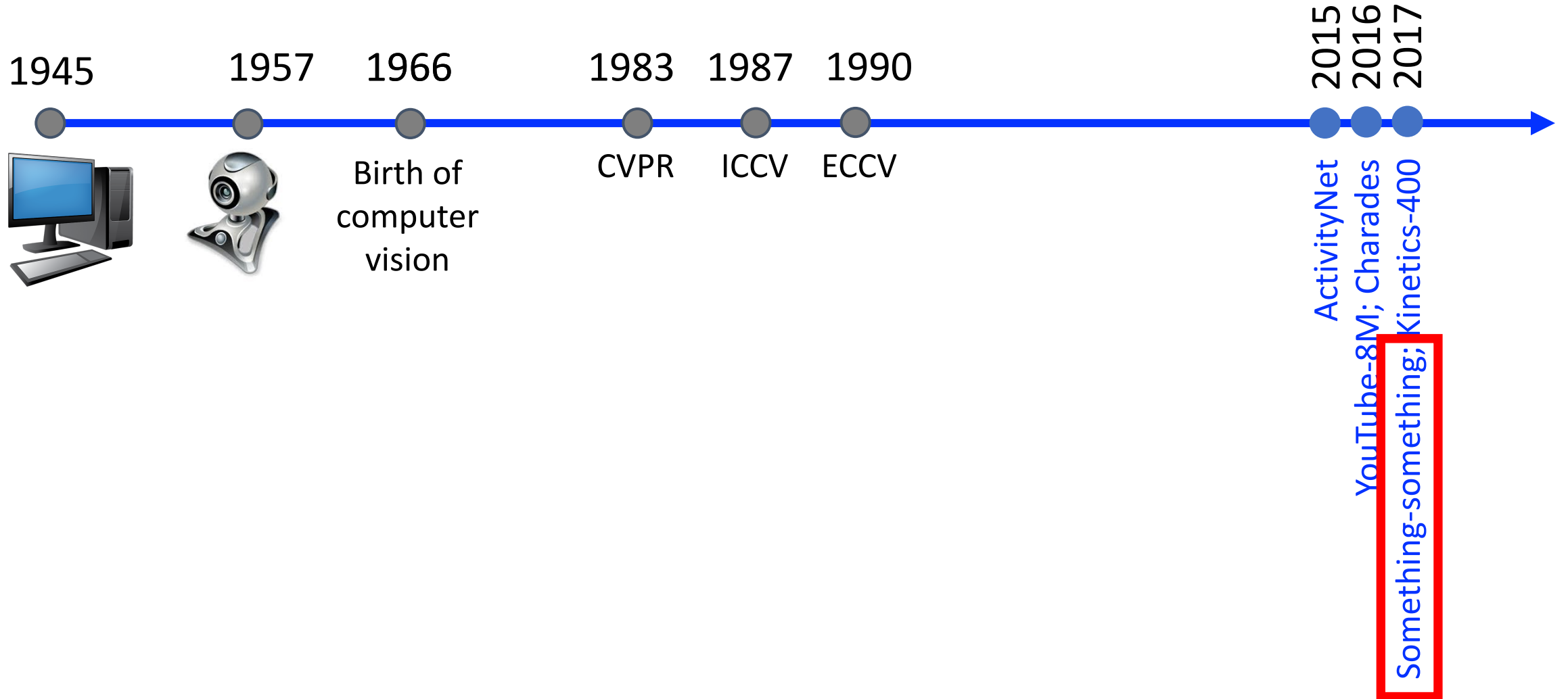
Recognize and locate activities taking place in a video



The Charades Activity Challenge aims towards automatic understanding of daily activities, by providing realistic videos of people doing everyday activities. [The Charades dataset](#) is collected for an unique insight into daily tasks such as drinking coffee, putting on shoes while sitting in a chair, or snuggling with a blanket on the couch while watching something

<http://vuchallenge.org/charades.html>

# Datasets



# Charades

Collection of videos to help models learn common sense features for predicting an activity label; (e.g., “opening” for blinds, door, mouth, zipper)

- More than 100,000 videos ranging from 2 to 6 seconds
- Represents 174 classes

# Something-something

## 1. Category Selection

\* Authors created 175 something-something templates

e.g.,

<b>10 selected classes</b>
Dropping [something]
Moving [something] from right to left
Moving [something] from left to right
Picking [something] up
Putting [something]
Poking [something]
Tearing [something]
Pouring [something]
Holding [something]
Showing [something] (almost no hand)



# Something-something

## 1. Category Selection

\* Authors created 175 something-something templates

## 2. Video Collection

\* Crowd workers submitted videos of them recording an implementation of the template

UI

You have selected 5 of 10 descriptions

- ▶ Folding something (2)
- ▶ Stuffing/Taking out (1)
- ▼ Holding something (5)
  - Holding [something]
  - Holding [something] over [something]
  - Holding [something] next to [something]
  - Holding [something] in front of [something]
  - Holding [something] behind [something]
- ▶ Crowd of things (2)
- ▶ Shadows (1)
- ▶ Collisions of objects (3)
- ▼ Tearing something (3)
  - Tearing [something] into two pieces
  - Pretending to be tearing [something that is not tearable]
  - Tearing [something] just a little bit
- ▶ Lifting/Tilting objects with other objects on them (3)
- ▼ Moving two objects relative to each other (4)
  - Moving [something] closer to [something]
  - Moving [something] away from [something]
  - Moving [something] and [something] closer to each other (fix the camera and use both hands to move both objects)
  - Moving [something] and [something] away from each other (fix the camera and use both hands to move both objects)
- ▶ Attaching/Trying to attach (2)
- ▶ Spinning something (3)
- ▶ Something falling (2)
- ▶ Putting/Taking objects into/out of/next to/... other objects (19)
- ▼ Rolling and sliding something (10)
  - Letting [something] roll down a slanted surface
  - Letting [something] roll up a slanted surface, so it rolls back down
  - Letting [something] roll along a flat surface
  - Putting [something] on a flat surface without letting it roll
  - Putting [something] that can't roll onto a slanted surface, so it stays where it is
  - Putting [something] that can't roll onto a slanted surface, so it slides down
  - Lifting a surface with [something] on it until it starts sliding down
  - Lifting a surface with [something] on it but not enough for it to slide down
  - Putting [something] onto a slanted surface but it doesn't glide down
  - Rolling [something] on a flat surface
- ▶ Plugging something into something (2)

You have uploaded 0 of 10 videos

Submit Task

holding something over something

Upload Video

moving something away from something

Upload Video

lifting a surface with something on it but not enough for it to slide down

Upload Video

moving something away from something

Upload Video

tearing something into two pieces

Upload Video

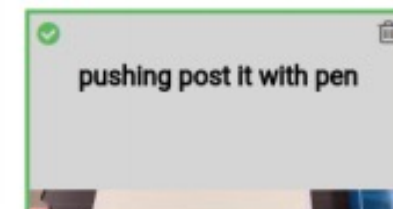
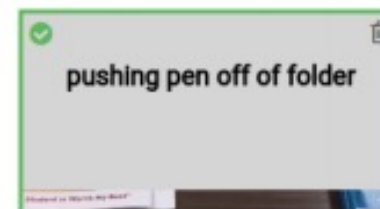
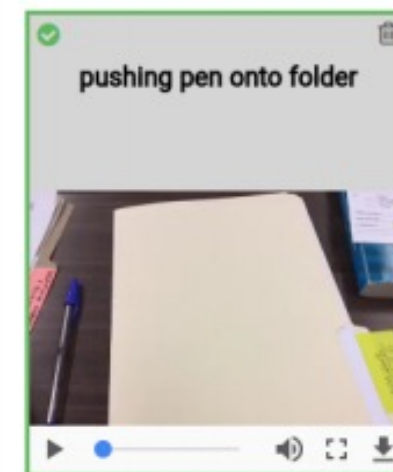
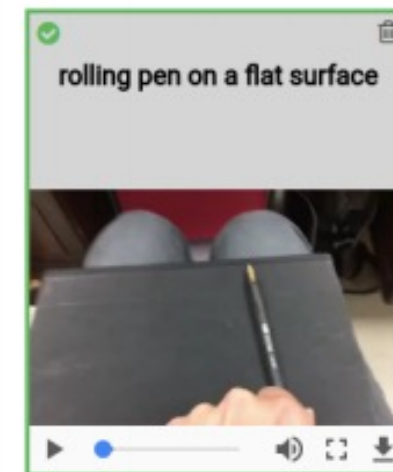
UI

You have selected 10 of 10 descriptions

- › Folding something (2)
- › Stuffing/Taking out (1)
- › Holding something (5)
- › Crowd of things (2)
- › Shadows (1)
- › Collisions of objects (3)
- › Tearing something (3)
- › Lifting/Tilting objects with other objects on them (3)
- › Moving two objects relative to each other (4)
- › Attaching/Trying to attach (2)
- › Spinning something (3)
- › Something falling (2)
- › Putting/Taking objects into/out of/next to/... other objects (19)
- › Rolling and sliding something (10)
- › Plugging something into something (2)
- › Twisting something (3)
- › Opening or closing something (4)
- › Pushing something (9)
- › Tipping something over (2)
- › Filming objects, without any actions (3)
- › Spilling something (3)
- › Turning something upside down (2)
- › Putting something somewhere (2)
- › Picking something up (2)
- › Hitting something with something (1)
- › Dropping something (5)
- › Poking something (9)
- › Throwing something (6)
- › Wiping something off of something (2)
- › Camera motions (8)
- › Showing objects and photos of objects (2)
- › Squeezing something (2)
- › Revolving something (2)

You have uploaded 10 of 10 videos

Submit Task



# Something-something Challenge

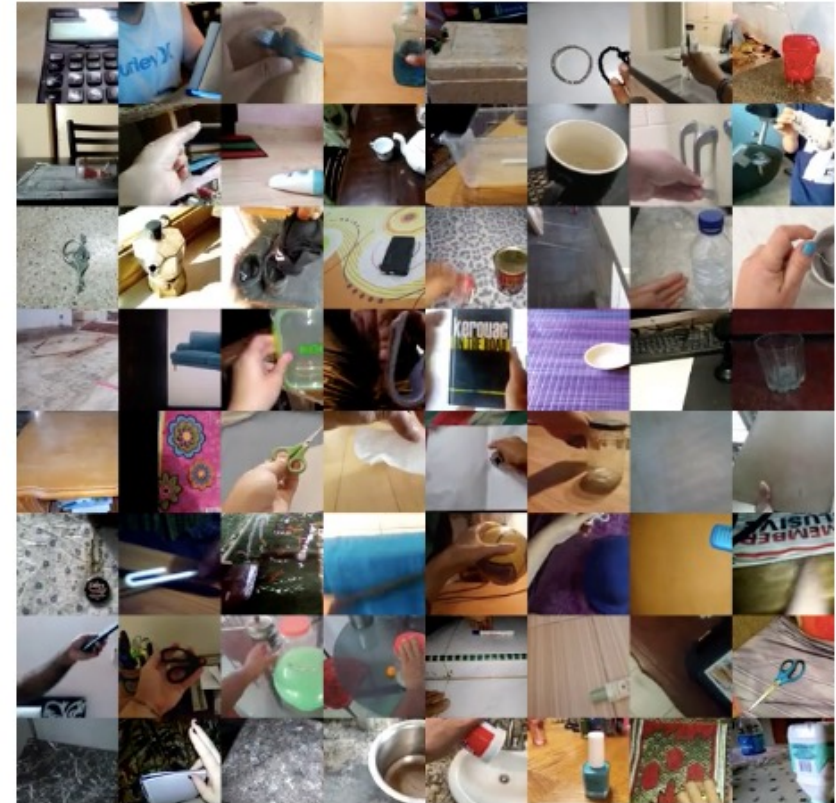
## The 20BN-something-something Dataset V2

### Introduction

The 20BN-SOMETHING-SOMETHING dataset is a large collection of densely-labeled video clips that show **humans performing pre-defined basic actions with everyday objects**. The dataset was created by a large number of crowd workers. It allows machine learning models to develop fine-grained understanding of basic actions that occur in the physical world. It is **available free of charge for academic research**. Commercial licenses are available upon request.

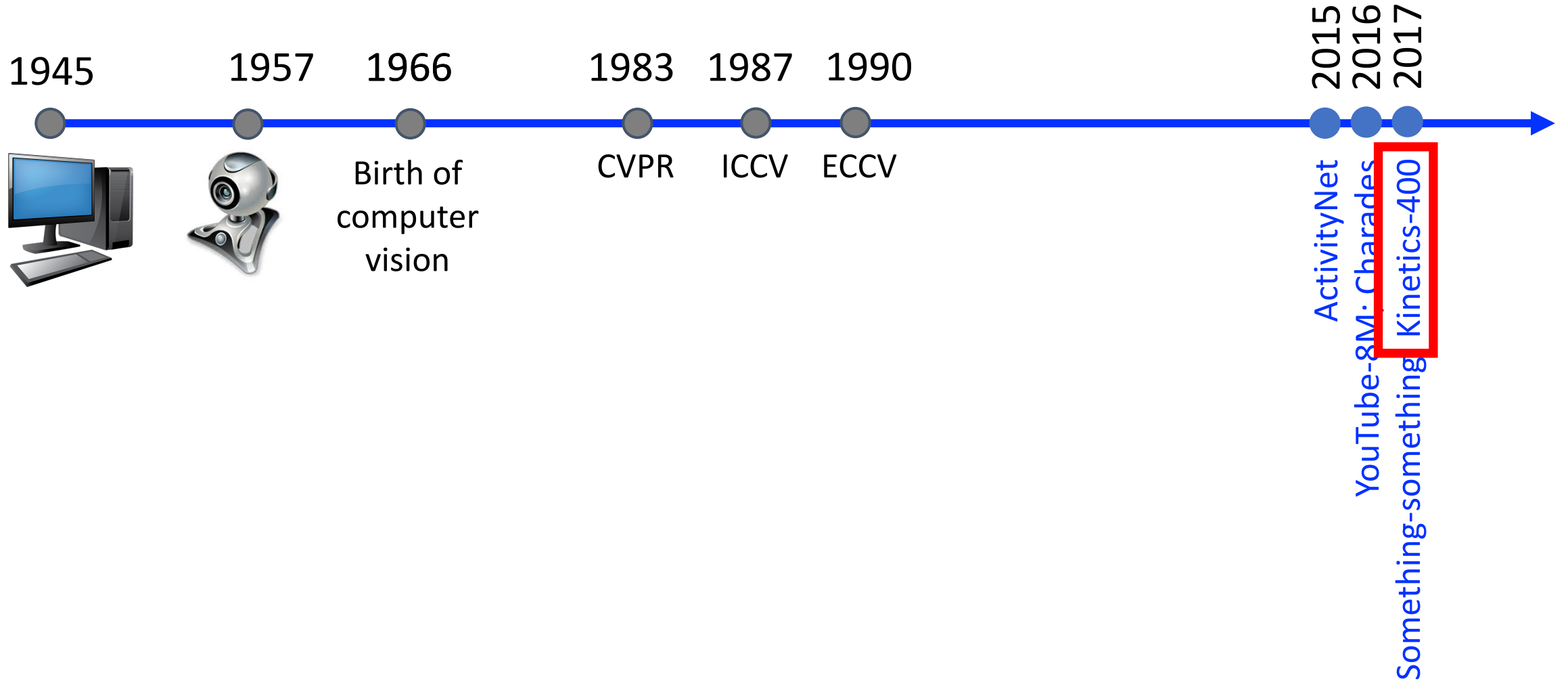
This is the second release of the dataset. The first release is also still available [here](#). The new release features the following updates:

- **Greatly increased number of videos:** With 220,847 videos (vs. 108,499 in V1) we release more than twice as many videos.
- **Object annotations and captioning:** For each video in the training and validation sets we now also provide object annotations in addition to the video label if applicable. For example, for a label like "Putting



<https://20bn.com/datasets/something-something>

# Datasets



# Kinetics-400

Multi-class classification dataset of human actions videos that is two orders of magnitude larger than prior work

- 306,245 videos that each roughly 10 seconds long
- Represents 400 classes covering with 400-1150 clips per class:
  - Person Actions; e.g., drawing, drinking, laughing, punching
  - Person-Person Actions; e.g., hugging, kissing, shaking hands;
  - Person-Object Actions; e.g., opening presents, mowing lawn, washing dishes

# Kinetics-400

## 1. Category Selection

\* Categories taken from existing datasets and AMT workers feedback about more suitable categories for clips with existing category labels assigned

## 2. Video Collection



\* Image action recognition classifiers applied to identify relevant clips (5 seconds before plus 5 seconds after image where activity is recognized)

## 3. Video Verification

\* AMT worker verifies action is present in each  
\* Label determined by majority of 5 workers

# Kinetics Challenge & Annual Workshop



[HOME](#)

[PEOPLE](#)

[CHALLENGE](#)

[PROGRAM](#)

[DATES](#)

[EVALUATION](#)

[CONTACT](#)



## Task B – Kinetics Challenge

[Challenge 2020](#) → **Task B – Kinetics Challenge**

The goal of the Kinetics dataset is to help the computer vision and machine learning communities advance models for video understanding. Given this large human action classification dataset, it may be possible to learn powerful video representations that transfer to different video tasks.

For information related to this task, please contact: [enoland@google.com](mailto:enoland@google.com), [joaluis@google.com](mailto:joaluis@google.com)

[http://activity-net.org/challenges/2020/tasks/guest\\_kinetics.html](http://activity-net.org/challenges/2020/tasks/guest_kinetics.html)



# Class Task: Video Classification Costs

Assume the task is to classify the presence of 10 activities in 1,000,000 30-second videos. **How much do you believe it will cost in US dollars to collect all the crowdsourced annotations for the datasets?**

# Video Classification: Today's Topics

- Problem
- Applications
- Datasets
- **Evaluation metric**
- Computer vision models

# Accuracy Metric

- Percentage of correct predictions

# Video Classification: Today's Topics

- Problem
- Applications
- Datasets
- Evaluation metric
- Computer vision models

# Key Idea

- Recall: a video is a series of images

157	153	174	168	150	152	129	151	172	161	155	156
155	182	163	74	75	62	33	17	110	210	180	154
180	180	50	14	34	6	10	33	48	106	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	68	137	251	237	239	239	228	227	87	71	201
172	105	207	233	233	214	220	239	238	98	74	206
157	153	174	168	150	152	129	151	172	161	155	156
180	154	10	168	134	11	31	62	22	148		
159	181	158	227	178	143	182	106	36	190		
55	180	236	231	149	178	228	43	95	234		
71	201	236	187	86	150	79	38	218	241		
74	206	227	210	127	102	36	101	255	224		
20	169	103	143	96	50	2	109	249	215		
22	148	1	81	47	0	6	217	255	211		
36	190	0	0	12	108	200	138	243	236		
95	234	177	121	123	200	175	13	96	218		
218	241										
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	86	150	79	38	218	241
190	224	147	108	227	210	127	102	36	101	255	224
190	214	173	66	103	143	96	50	2	109	249	215
187	196	236	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

Time 1

1 minute

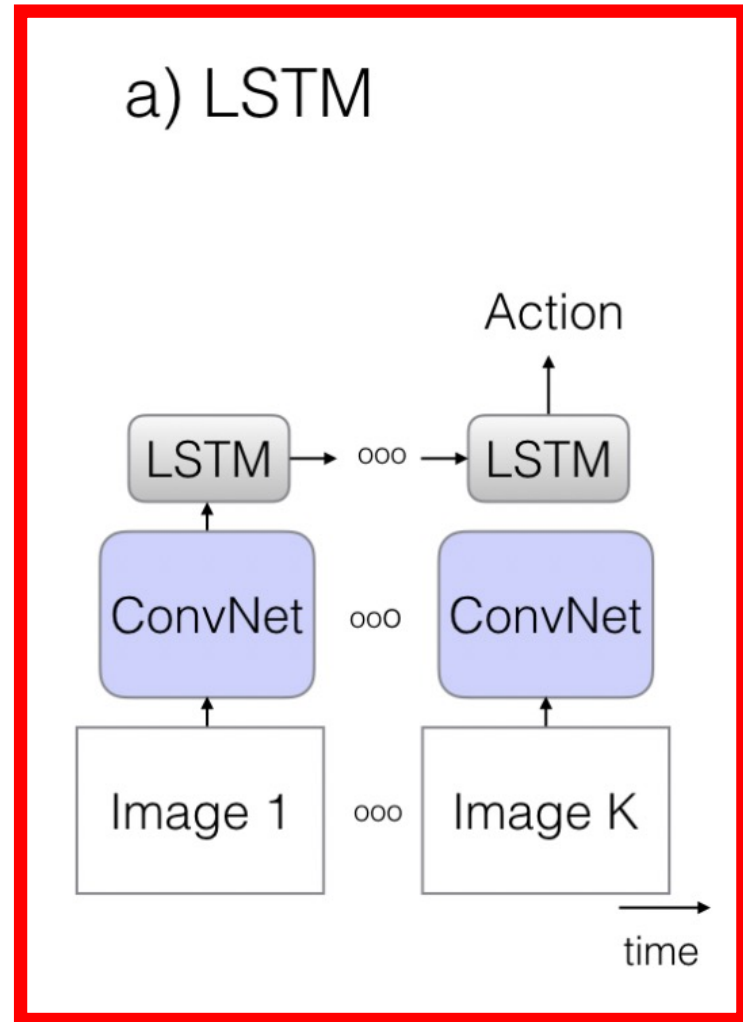
Analogous to:



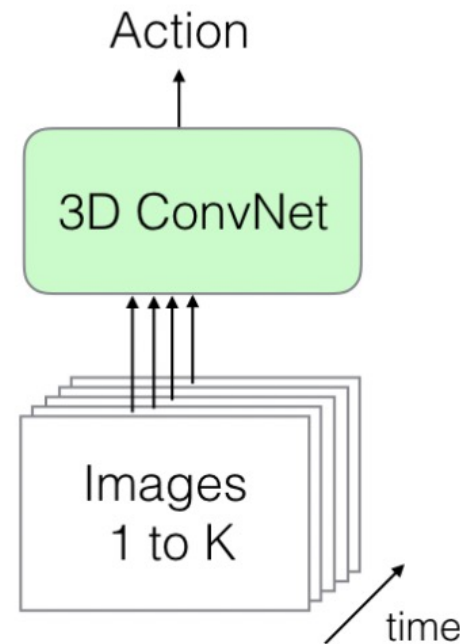
- How to go beyond image-based techniques and improve techniques by **considering the temporal relationship between frames in a video?** (e.g., recognize difference between a door that is opening vs closing)

# Approaches to Capture Temporal Information

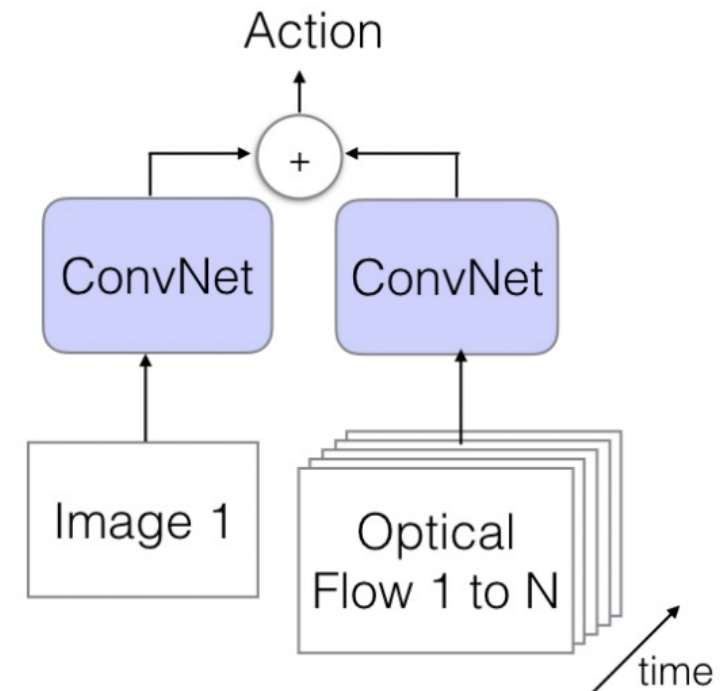
LSTM is a type of recurrent neural network (rnn); more on this Wednesday!



b) 3D-ConvNet



c) Two-Stream

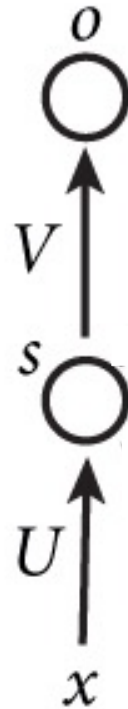


# Recurrent Neural Networks (RNNs)

- Main idea: use hidden state to **capture information about the past**

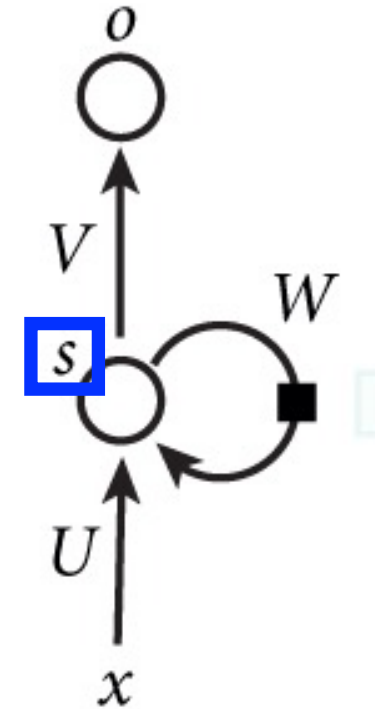
## Feedforward Network

Each layer receives input from the previous layer with no loops



## Recurrent Network

Each layer receives input from the previous layer and the output from the previous time step



# Recurrent Neural Networks (RNNs)

- Main idea: use hidden state to **capture information about the past**

Recurrence formula applied at every time step:

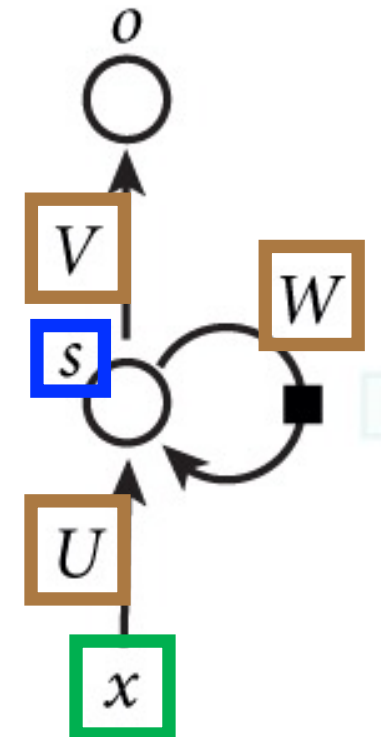
Model parameters

$$s_t = f_m(s_{t-1}, x_t)$$

New state                  Old state    Input at time step

## Recurrent Network

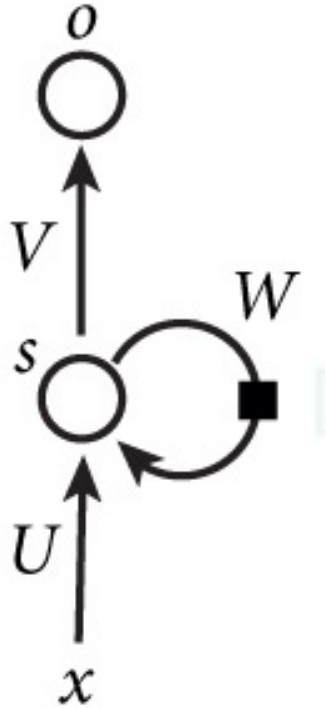
Each layer receives input from the previous layer and the output from the previous time step





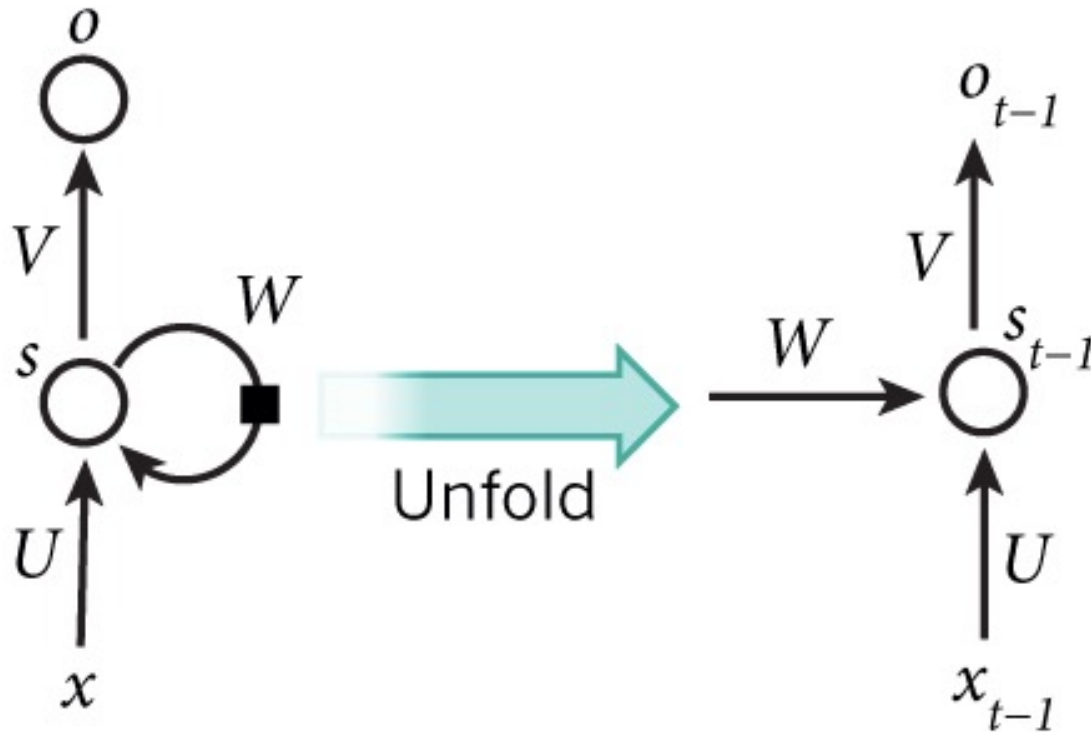
# RNN: Time Step 1

- Main idea: use hidden state to capture information about the past



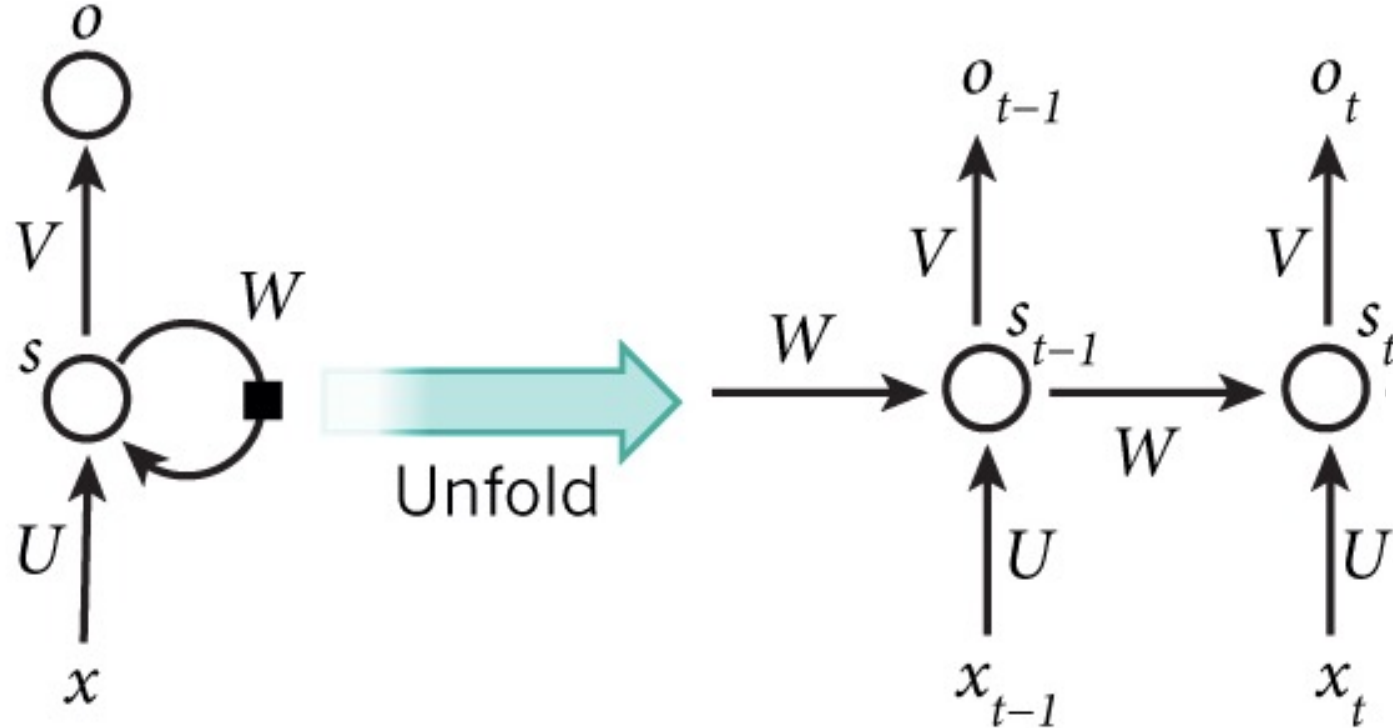
# RNN: Time Step 1

- Main idea: use hidden state to capture information about the past



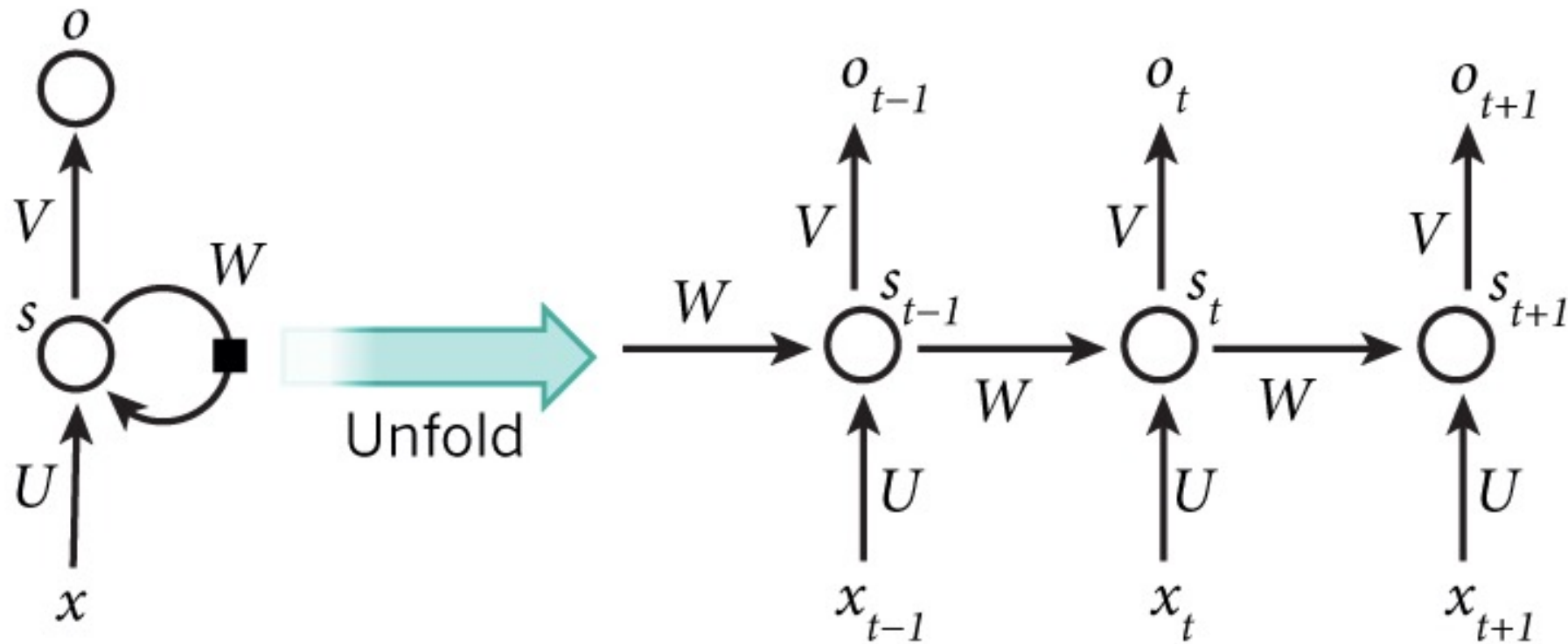
# RNN: Time Step 2

- Main idea: use hidden state to capture information about the past



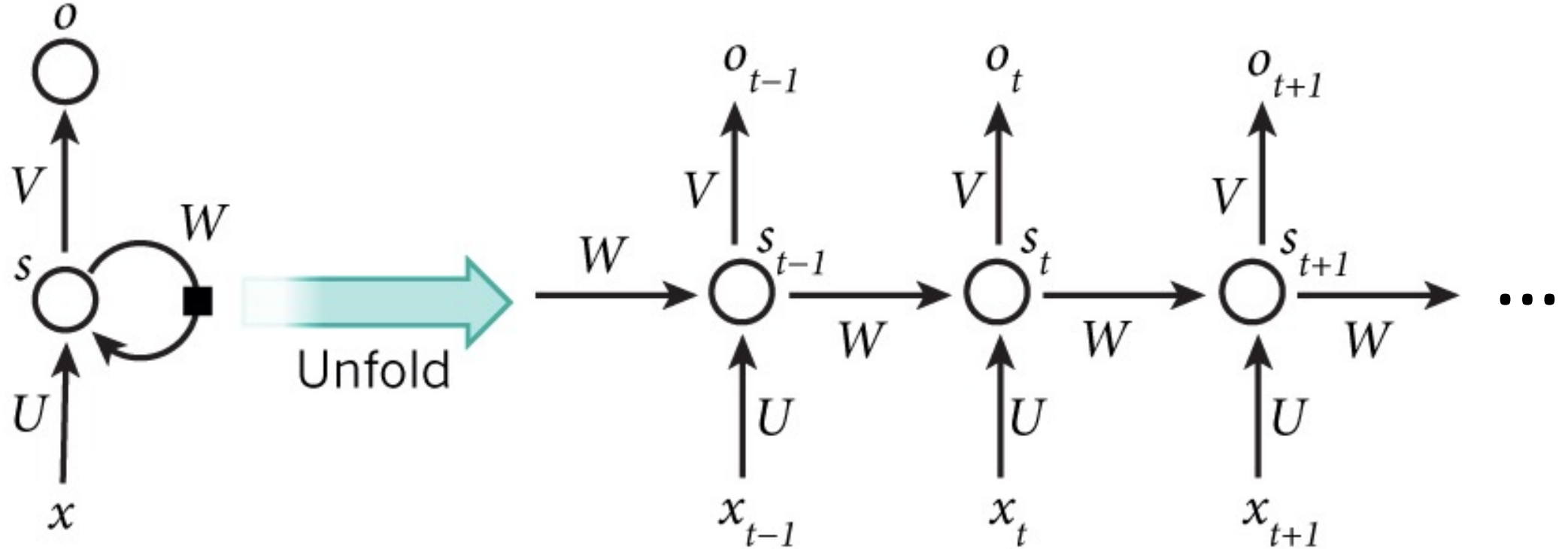
# RNN: Time Step 3

- Main idea: use hidden state to capture information about the past



# RNN: And So On...

- Main idea: use hidden state to capture information about the past



# RNN: Model Parameters and Inputs

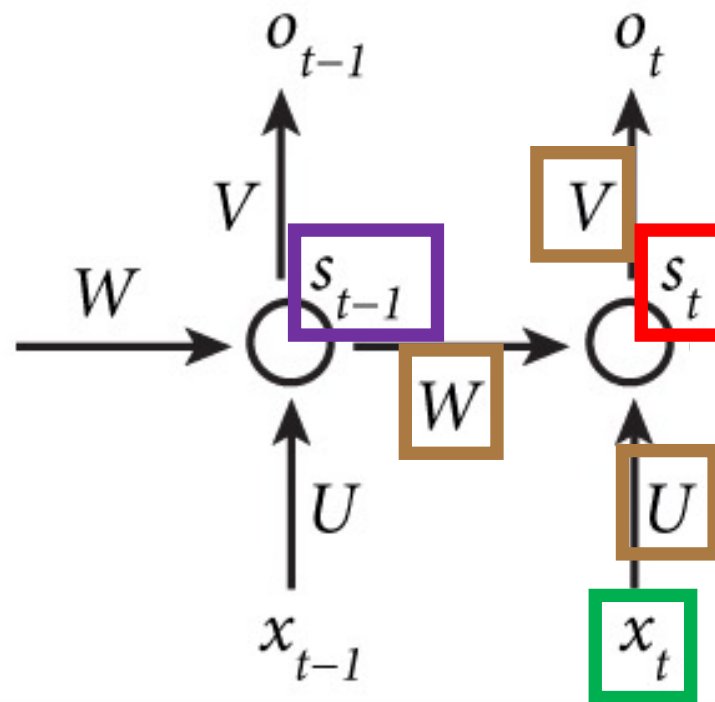
- Main idea: use hidden state to capture information about the past

Recurrence formula applied at every time step:

Model parameters

$$s_t = f_m(s_{t-1}, x_t)$$

New state                  Old state    Input at time step



# RNN: Model Parameters and Inputs

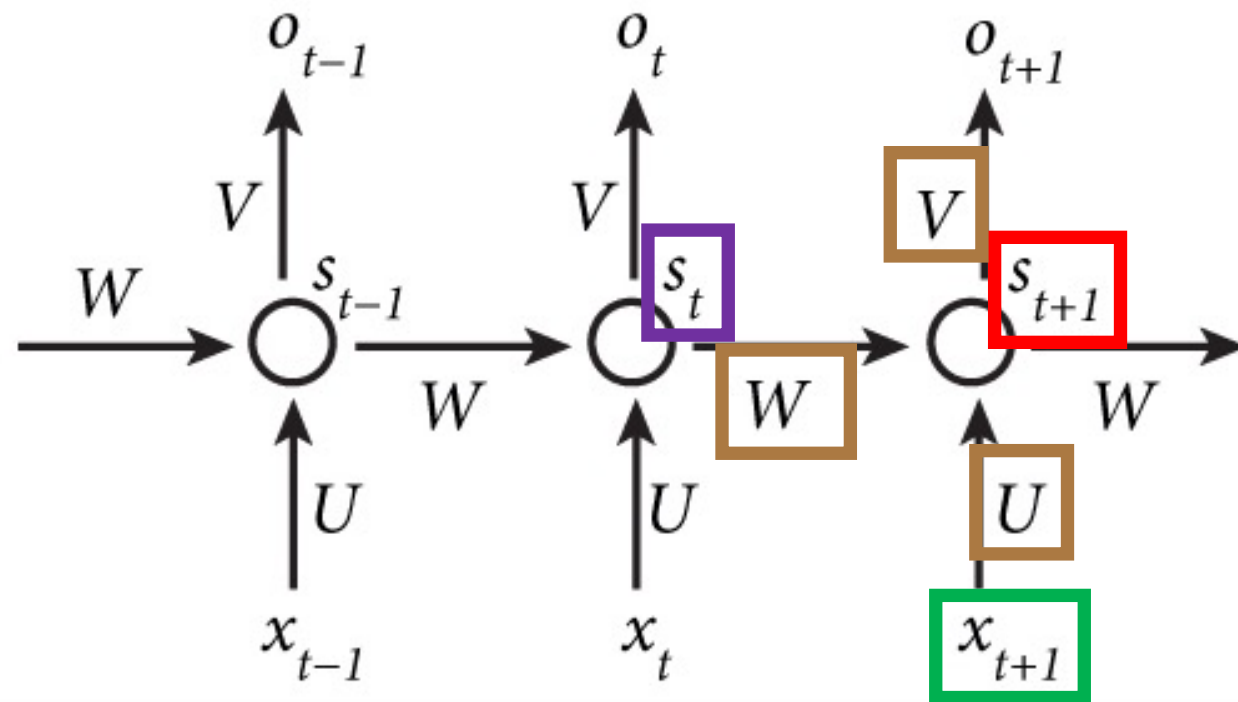
- Main idea: use hidden state to capture information about the past

Recurrence formula applied at every time step:

Model parameters

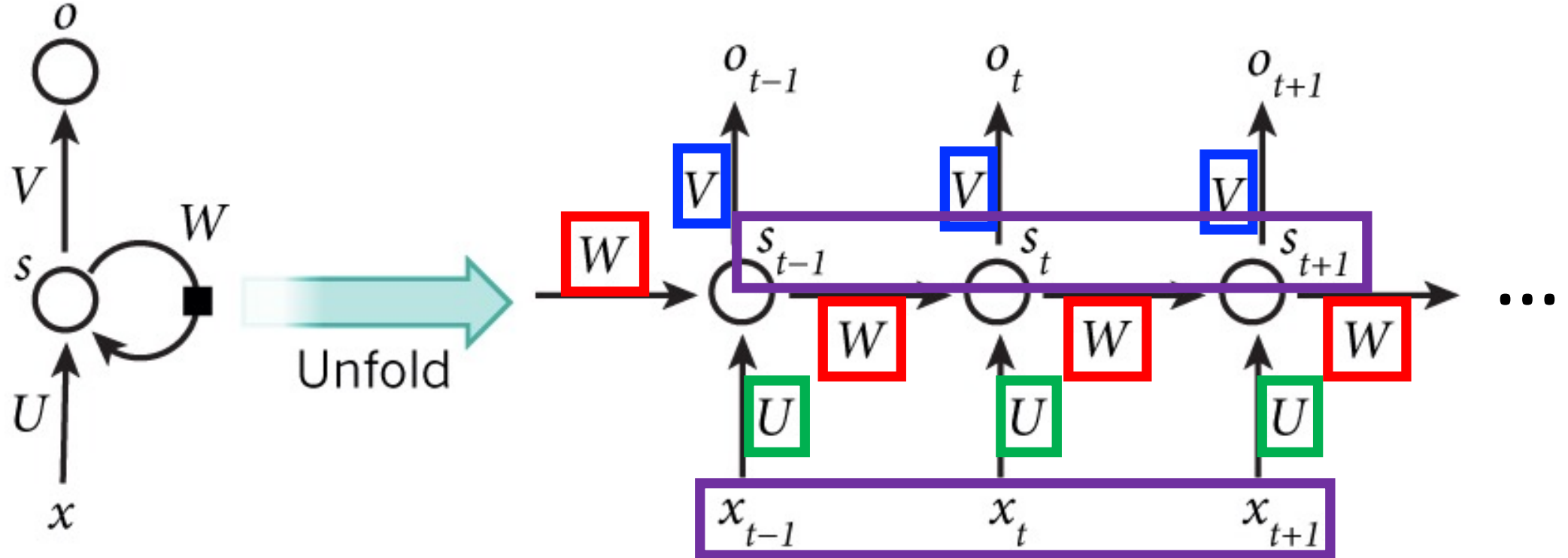
$$s_t = f_m(s_{t-1}, x_t)$$

New state                  Old state    Input at time step



# RNN: Model Parameters and Inputs

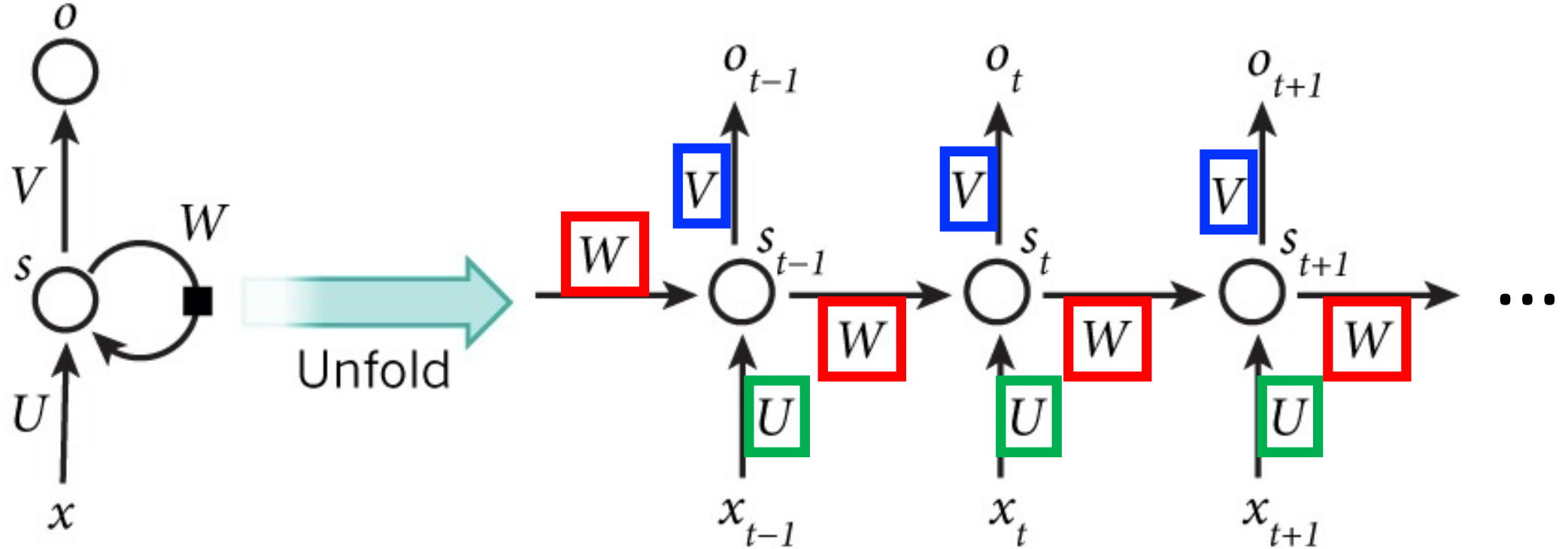
- All layers share the same model parameters ( $U$ ,  $V$ ,  $W$ )
  - What is different between the layers?





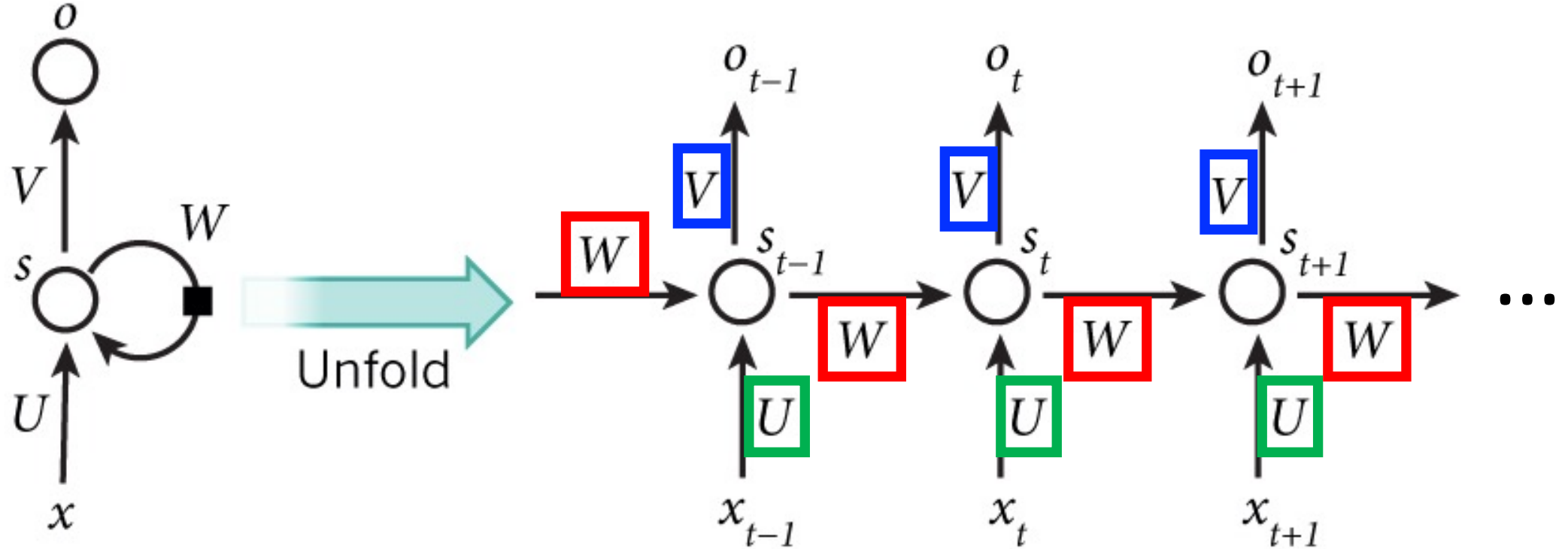
# RNN: Model Parameters and Inputs

- When unfolded, a RNN is a deep feedforward network with shared weights!

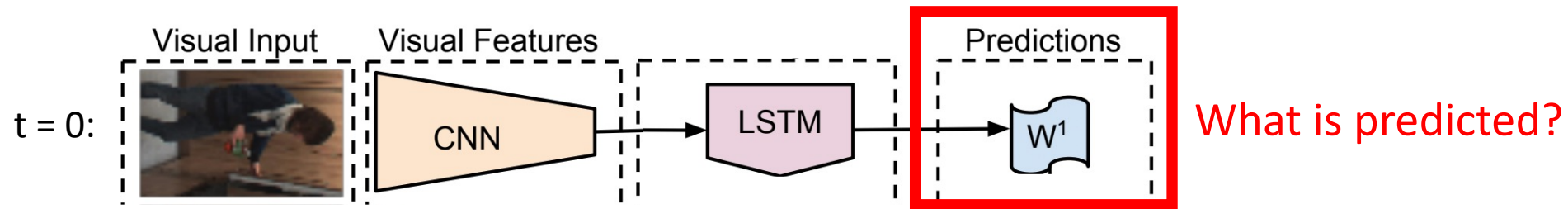


# RNN: Advantage

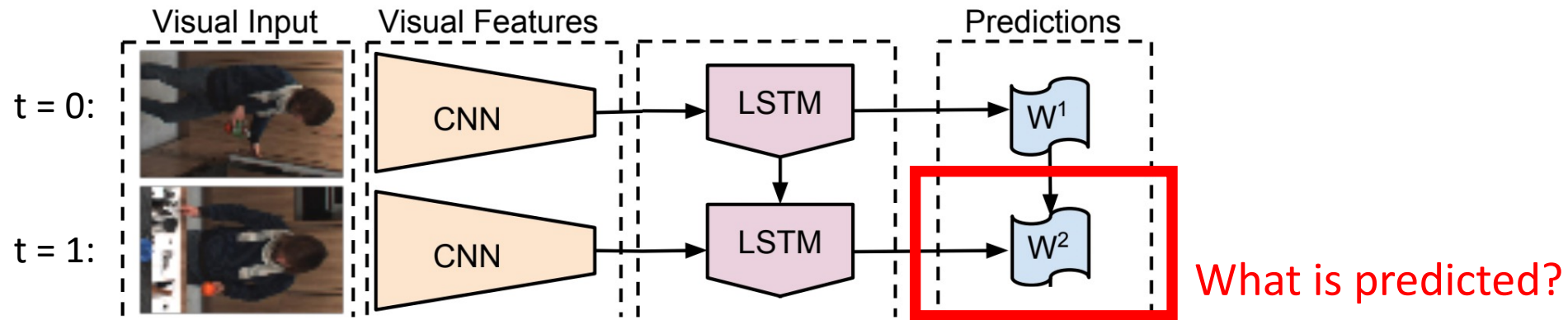
- Retains information about past inputs for an amount of time that depends on the model's weights and input data rather than a fixed duration selected a priori



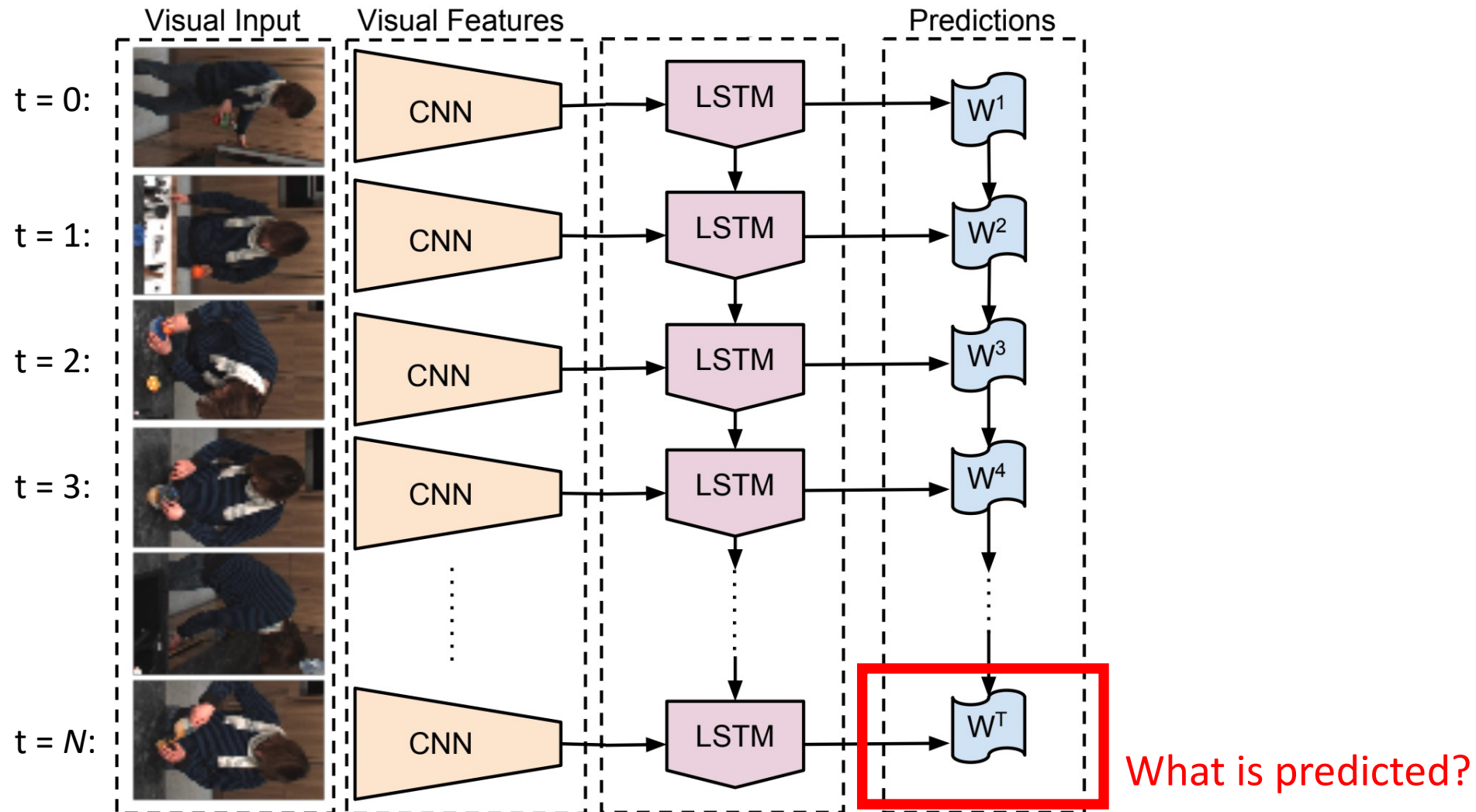
# RNN for Video Classification



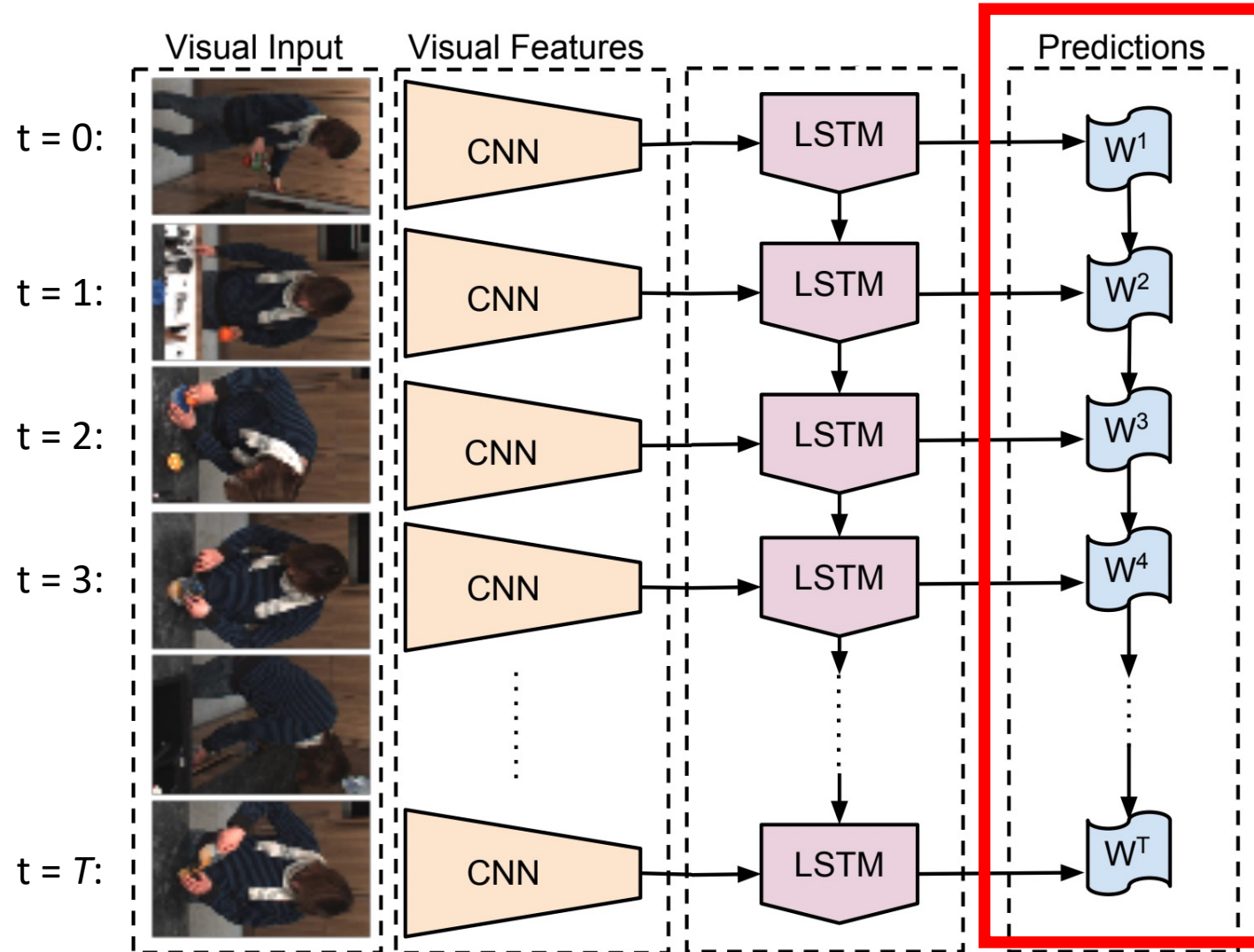
# RNN for Video Classification



# RNN for Video Classification

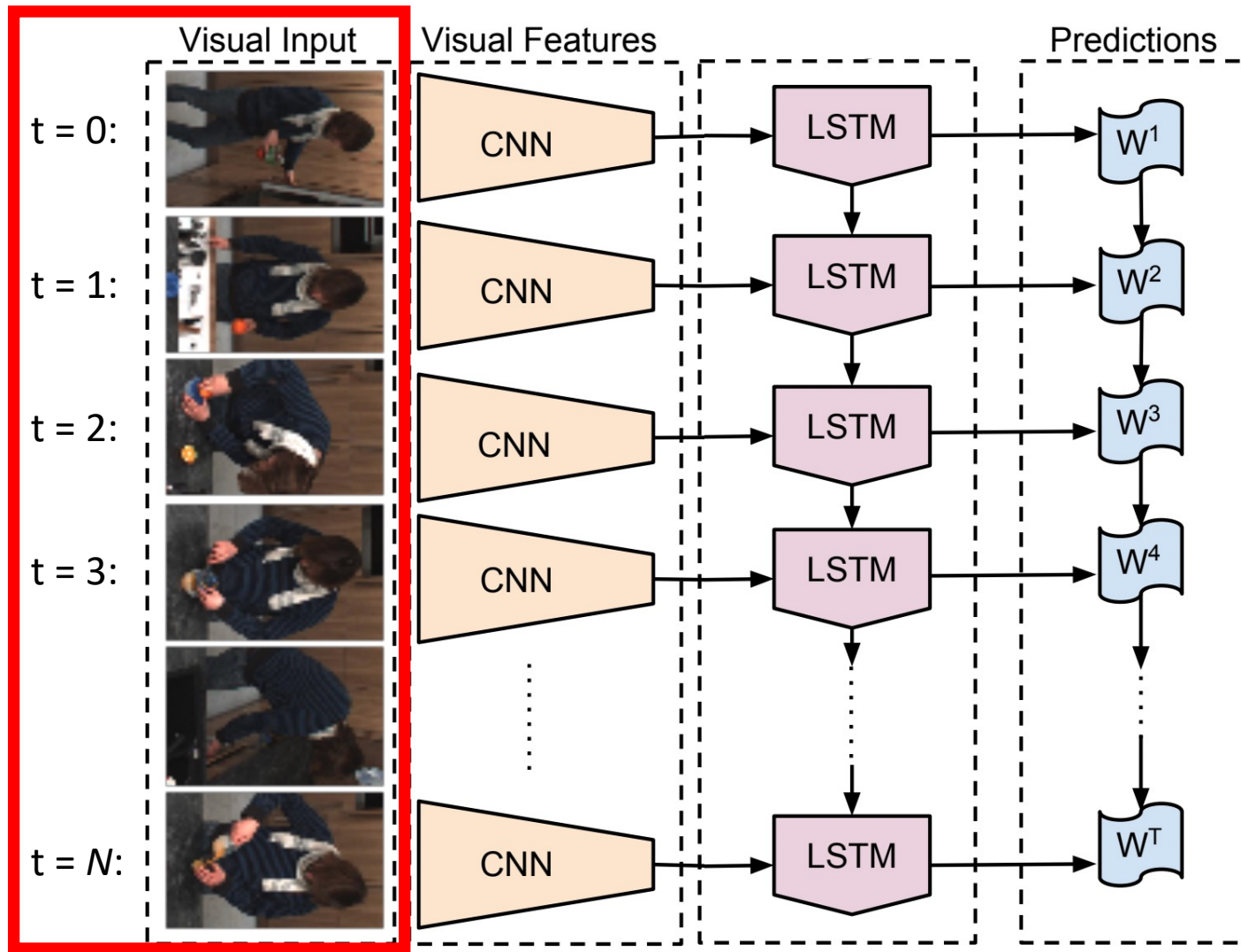


# RNN for Video Classification



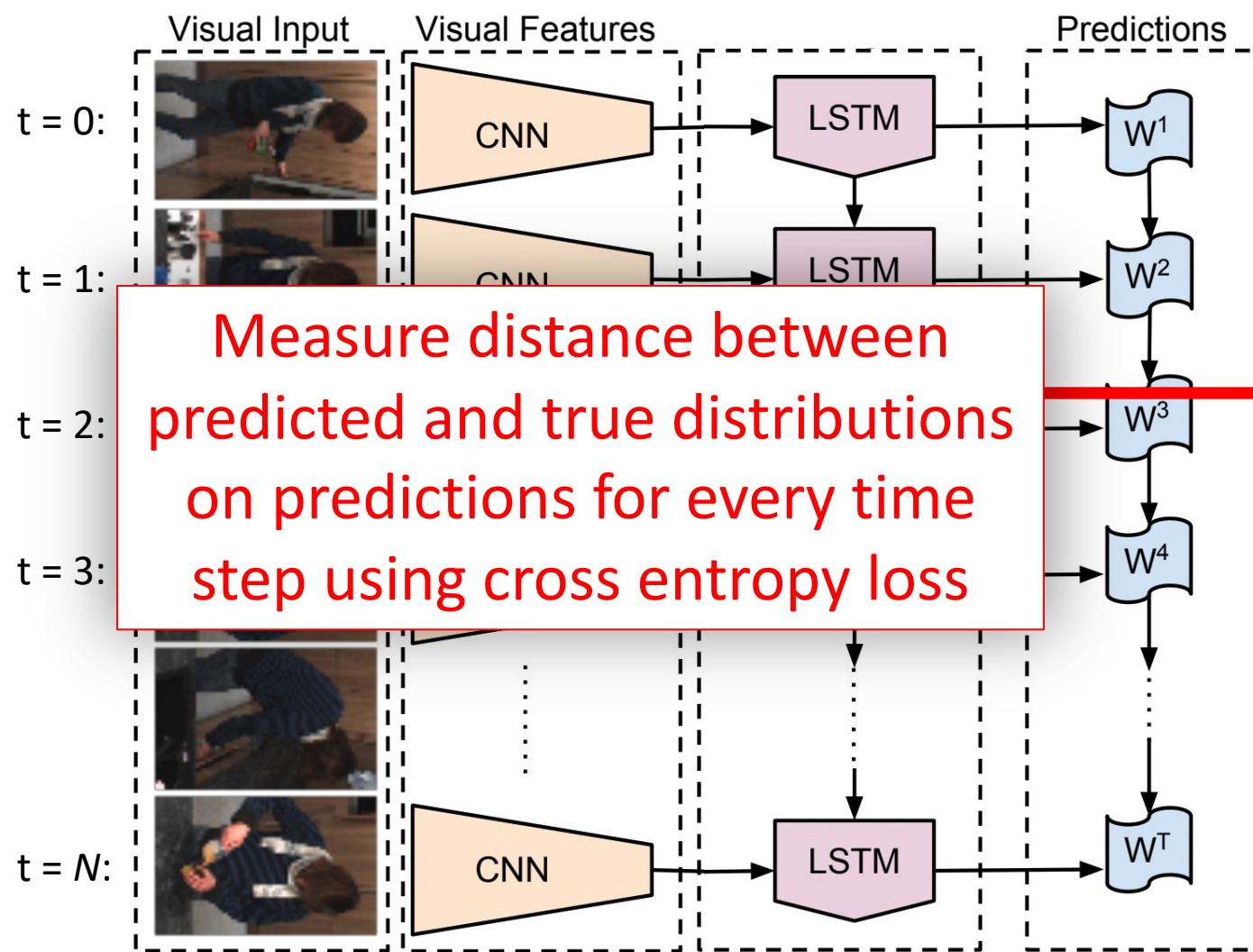
How should the final classification prediction be made (recall predictions are made from time step 1 to time step T)?

# RNN for Video Classification



What input duration is supported?

# RNN for Video Classification: Training Algorithm



- Repeat until stopping criterion met:

1. **Forward pass:** propagate training data through model to make prediction

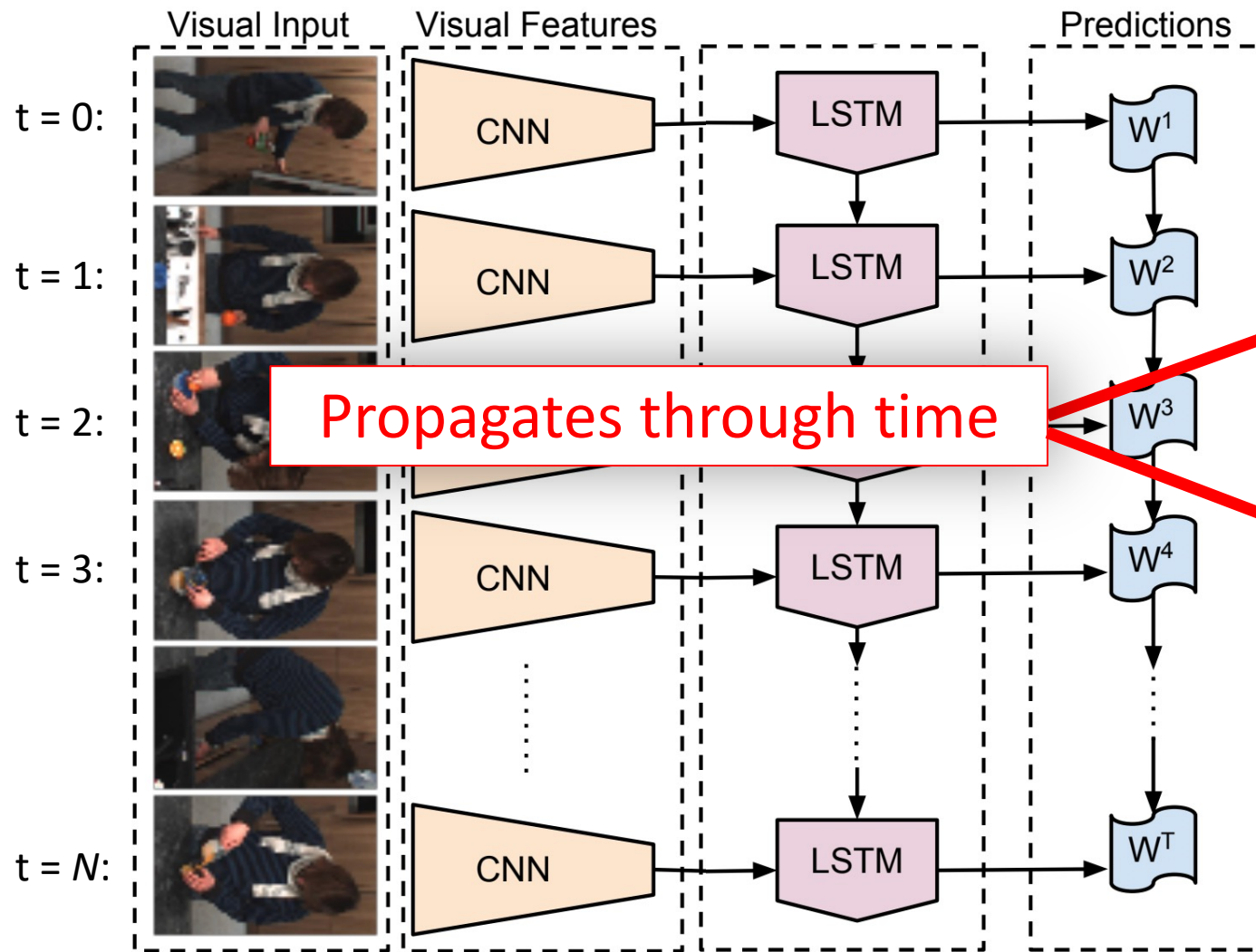
2. **Quantify the dissatisfaction with a model's results on the training data**

3. **Backward pass:** using predicted output, calculate gradients backward to assign blame to each model parameter

4. Update each parameter using calculated gradients



# RNN for Video Classification: Training Algorithm



- Repeat until stopping criterion met:

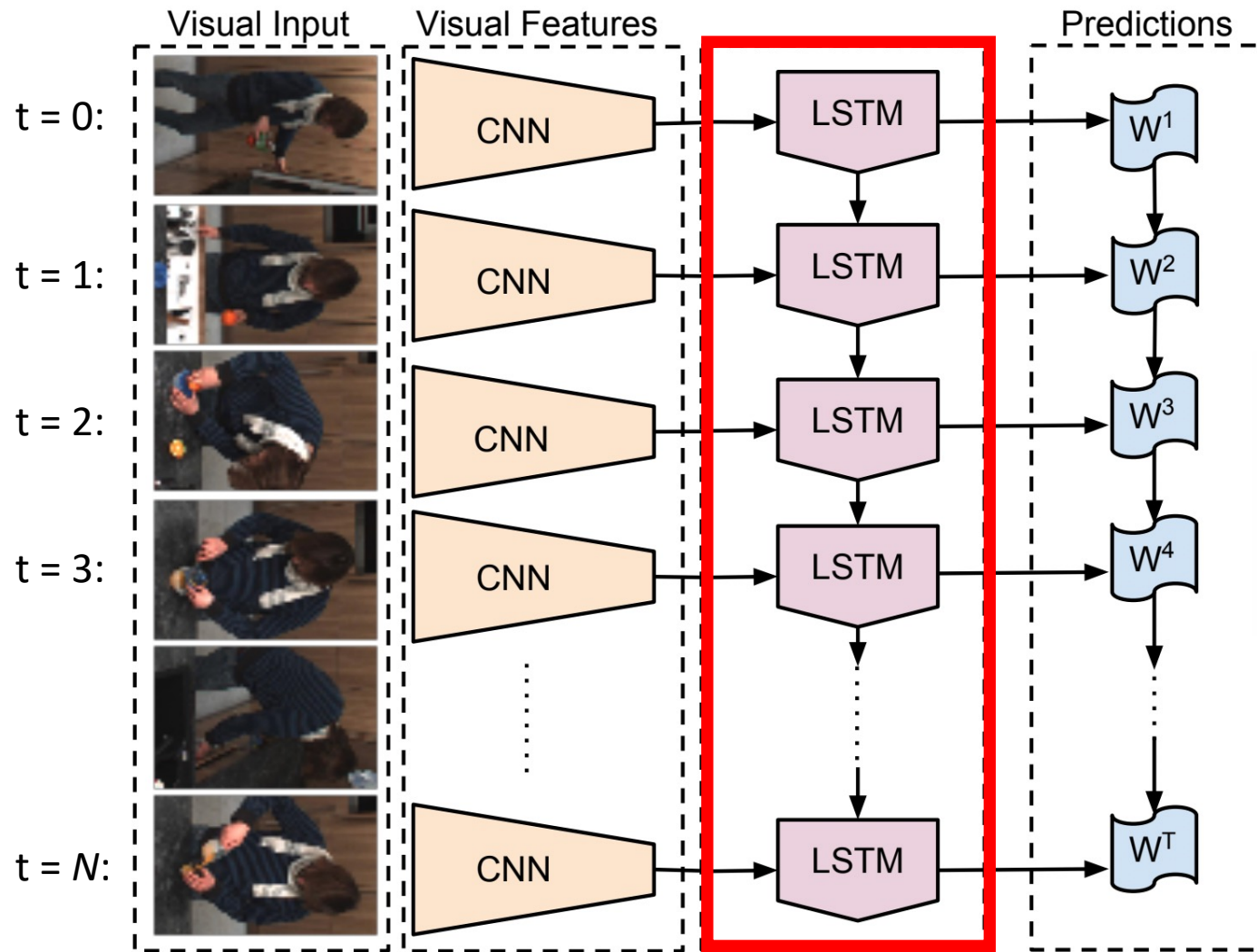
1. **Forward pass:** propagate training data through model to make prediction

2. Quantify the dissatisfaction with a model's results on the training data

3. **Backward pass:** using predicted output, calculate gradients backward to assign blame to each model parameter

4. Update each parameter using calculated gradients

# RNN for Video Classification

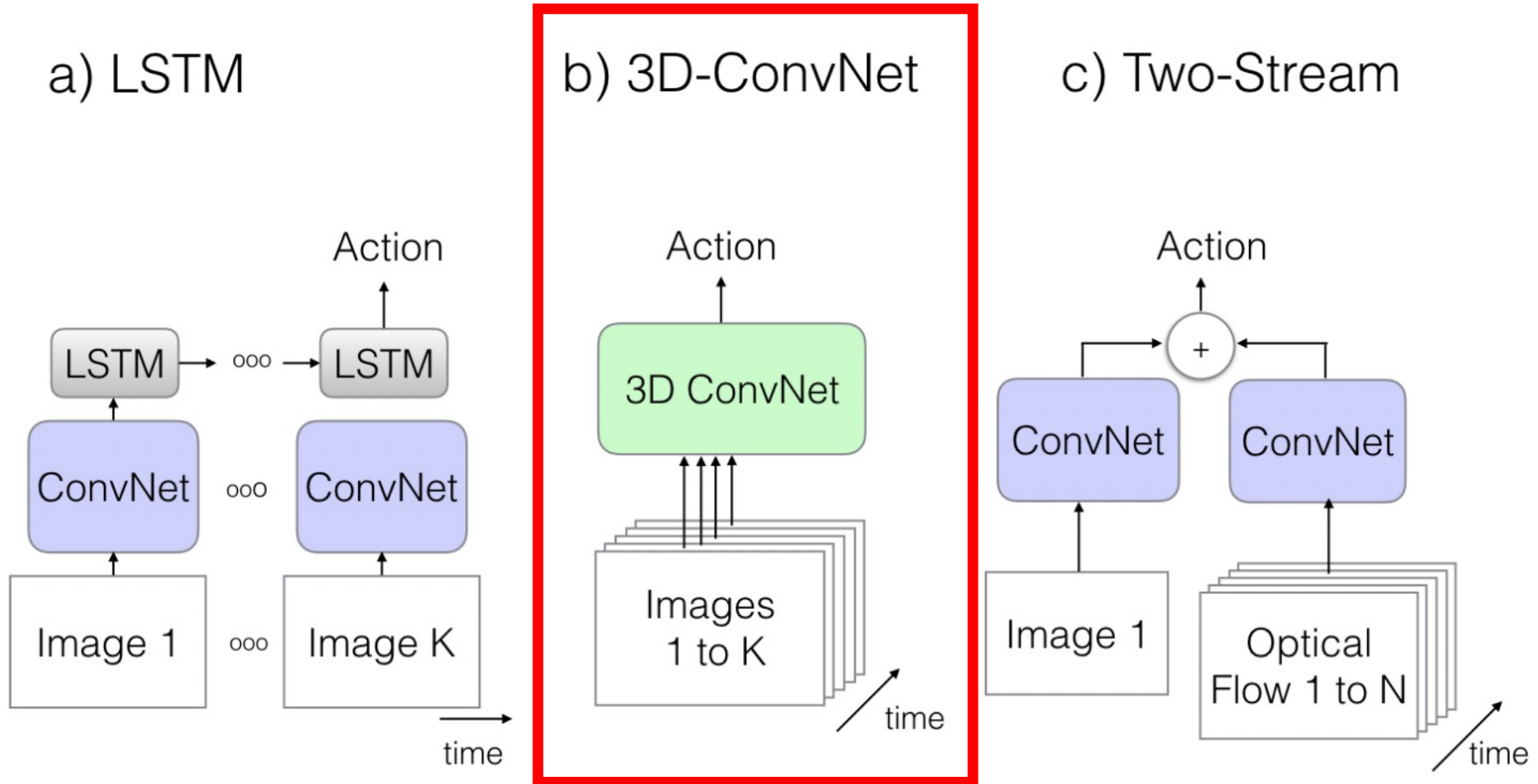


Recall: The LSTM layer's weights and input data determine what information about the past gets propagated to later time steps

# RNN for Video Classification: Limitations

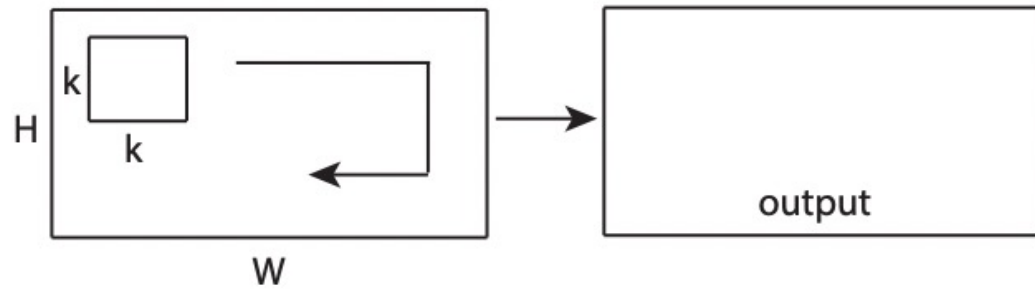
- Successful training requires many videos which makes RNNs resource-hungry and time-consuming

# Approaches to Capture Temporal Information

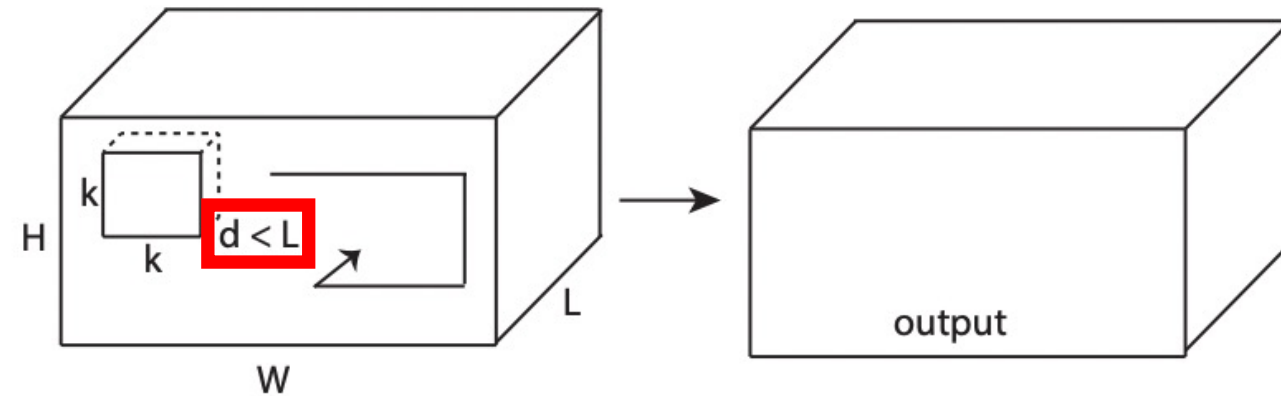


# ConvNet Architecture

- CNNs with **3D kernels** instead of 2D kernels to preserve temporal information in addition to spatial information



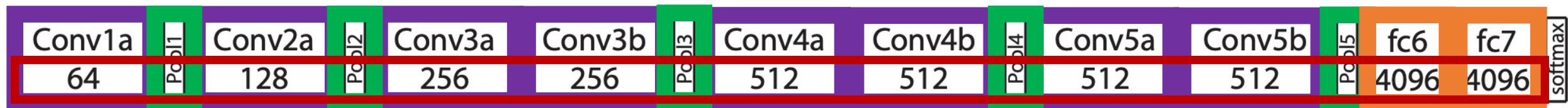
(a) 2D convolution



(c) 3D convolution

# ConvNet Architecture

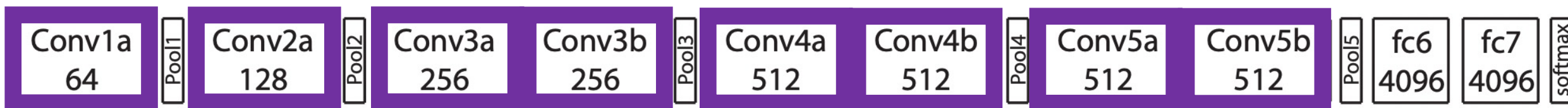
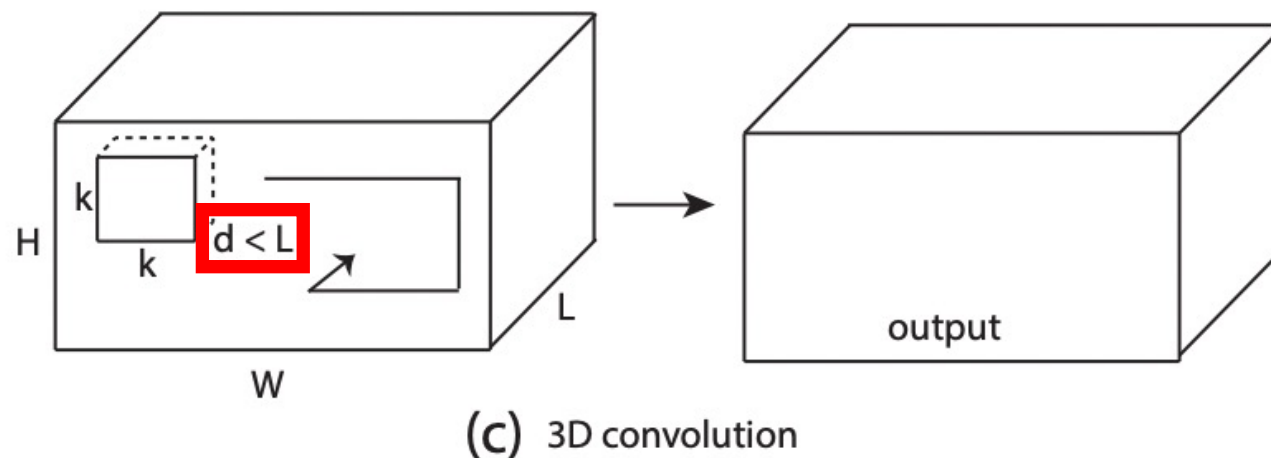
- How many **convolutional layers** are there?
- How many **pooling layers** are there?
- How many **fully-connected layers** are there?



Numbers indicate number of kernels/nodes per layer

# ConvNet Architecture

- Key question: what kernel depth should be used?
  - From experimentation: 3 at every layer (i.e., 3x3x3 kernels)



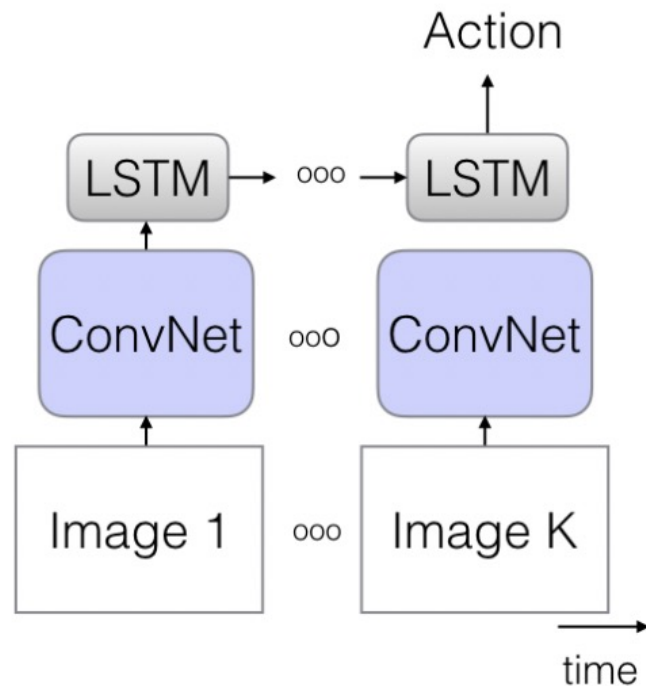
# ConvNet: Limitations

- 3D kernels introduce many model parameters and so successful training requires many videos (i.e., resource-hungry and time-consuming)
- Does not capture long-term temporal information (duration determined by depth of kernel)

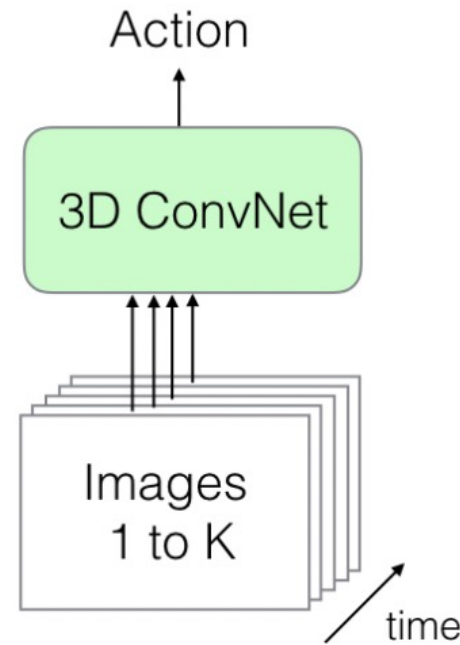


# Approaches to Capture Temporal Information

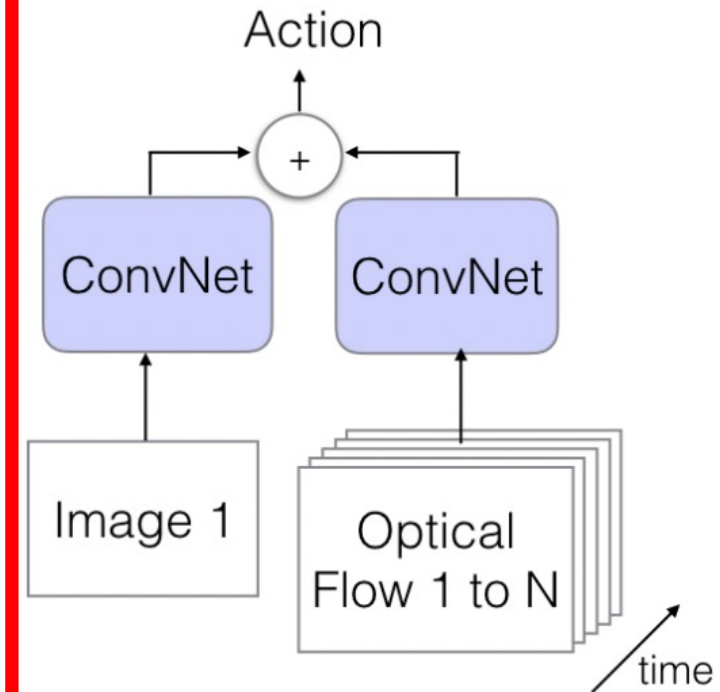
a) LSTM



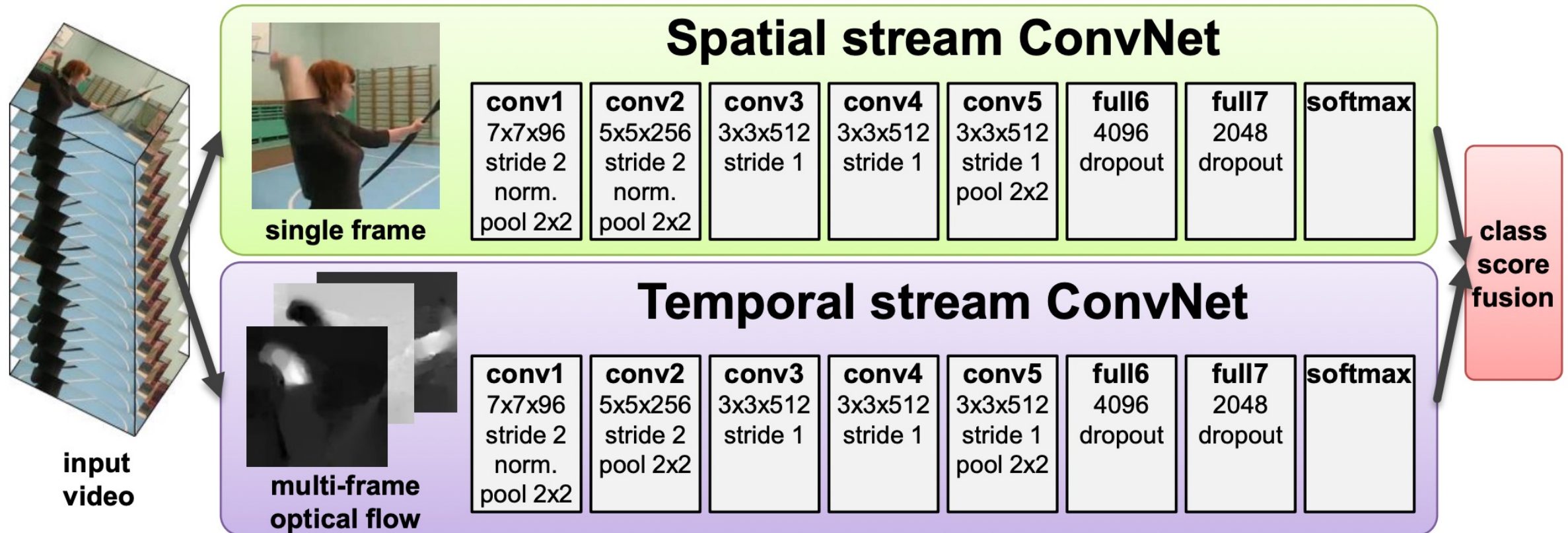
b) 3D-ConvNet



c) Two-Stream

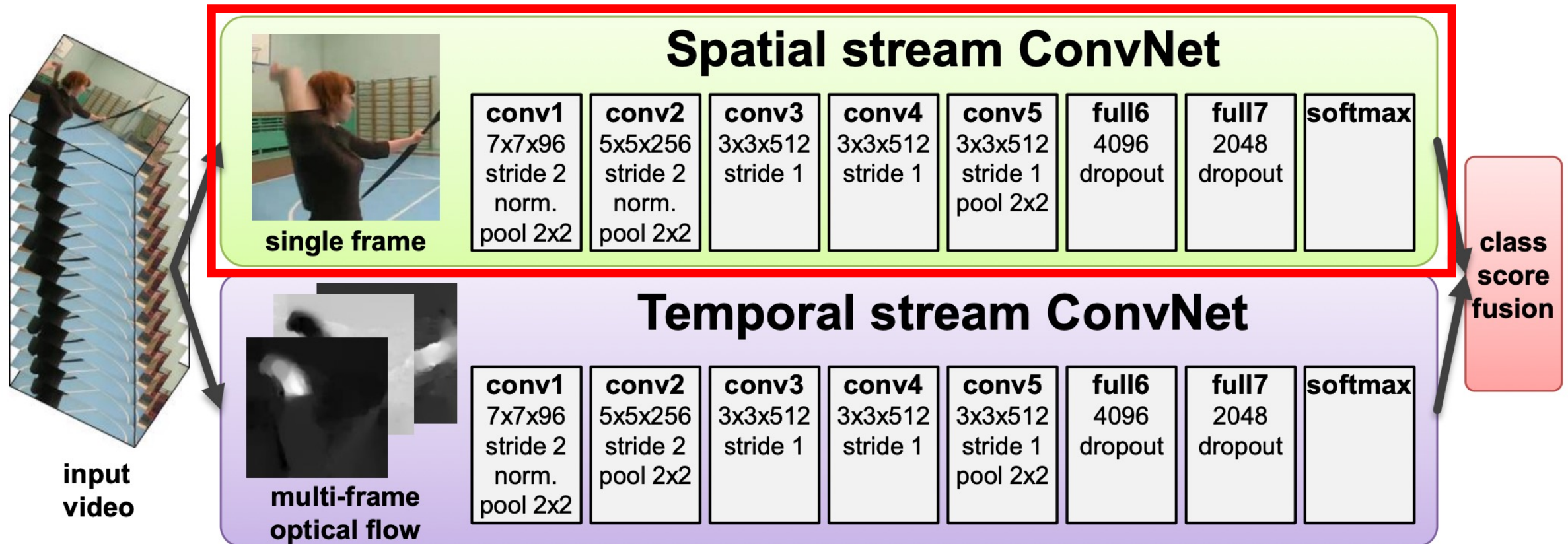


# Two-Stream Architecture



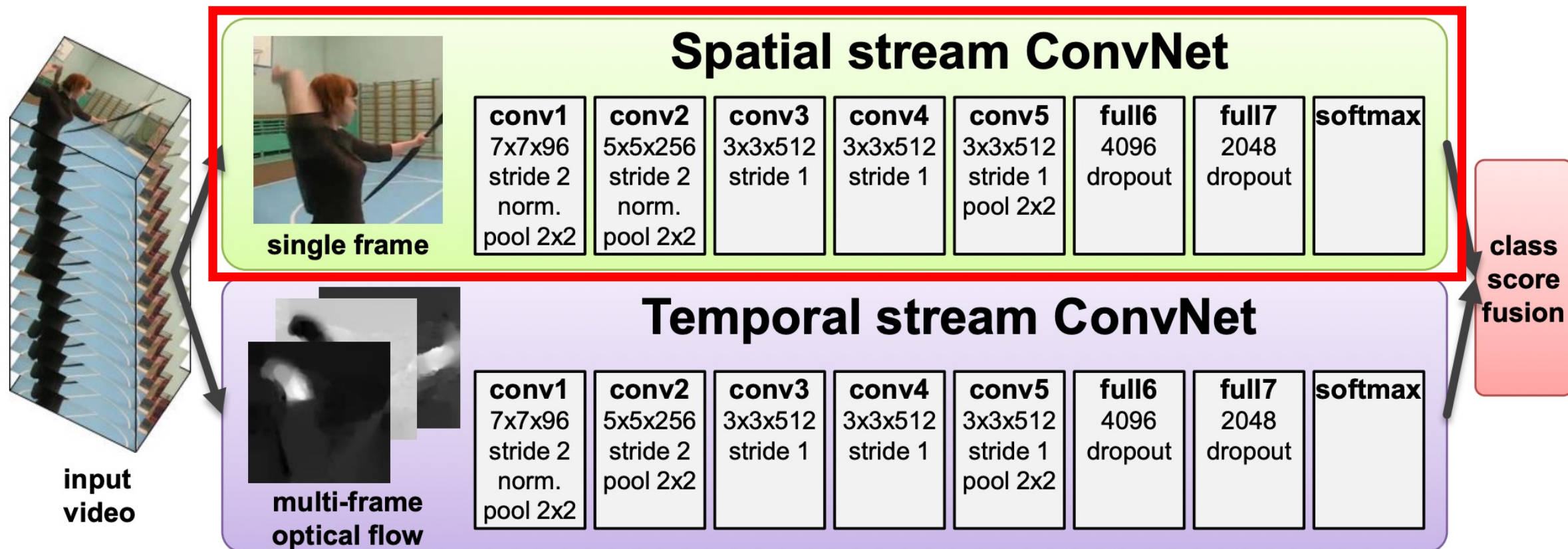
# Two-Stream Architecture

Learns to predict from still images the actions

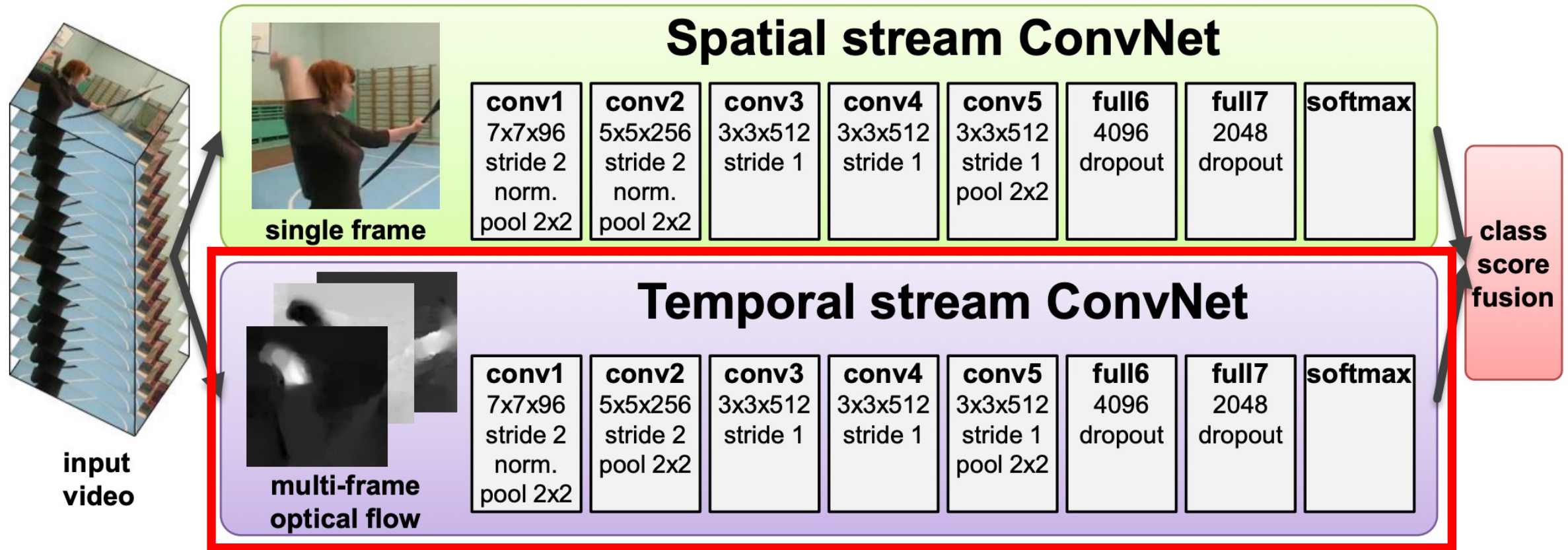


# Two-Stream Architecture

Pre-trained on ImageNet

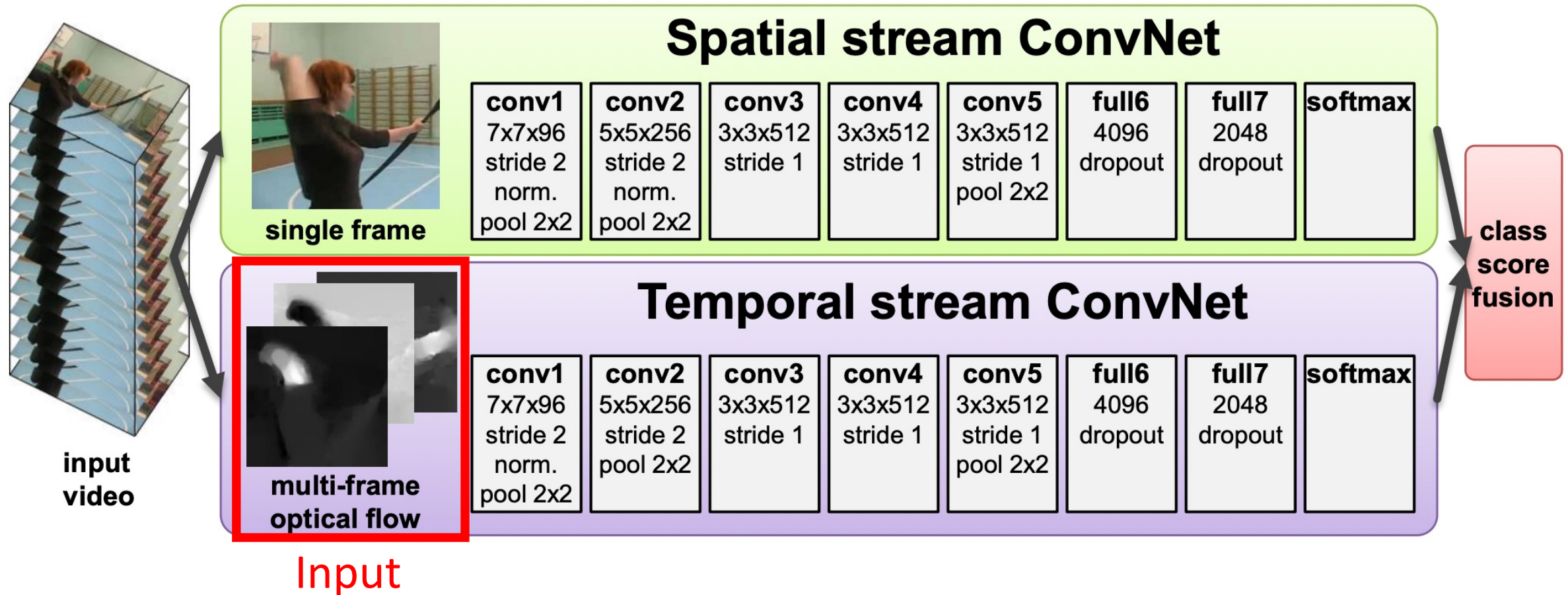


# Two-Stream Architecture

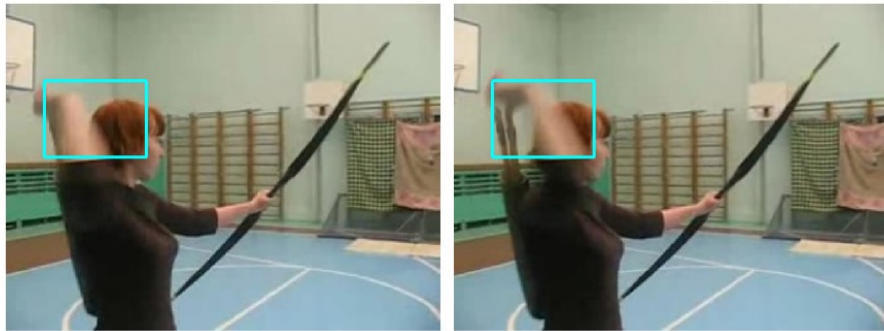


Learns to predict from explicit motion representations the actions

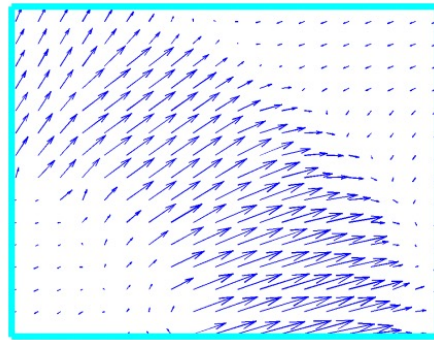
# Two-Stream Architecture



# Two-Stream Architecture: Input (Optical Flow)

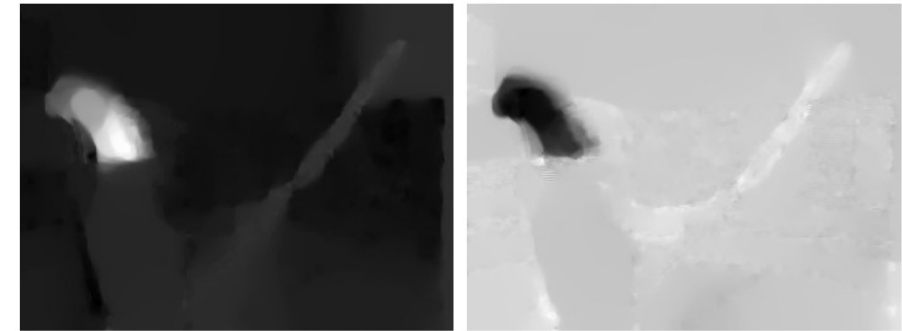


Two consecutive frames



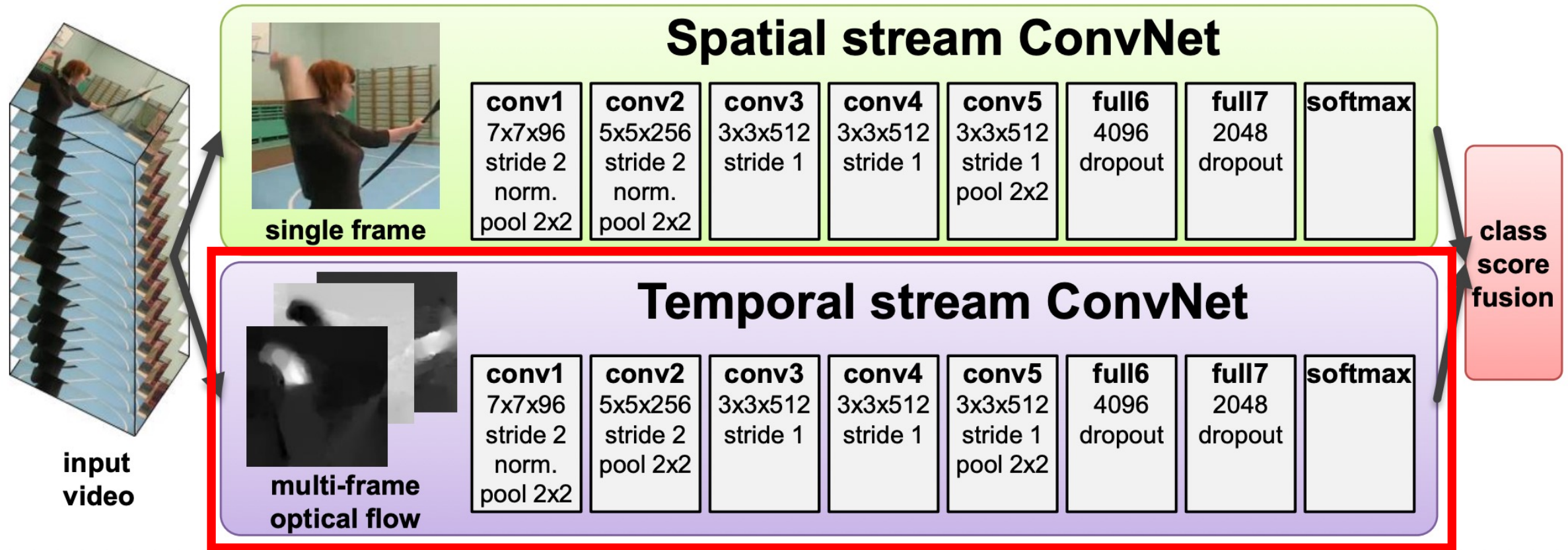
Vector fields showing where each point in the original frame moved in the subsequent frame

Input: stack pairs from consecutive frames



Vector fields decomposed into their horizontal (left) and vertical components (right)

# Two-Stream Architecture

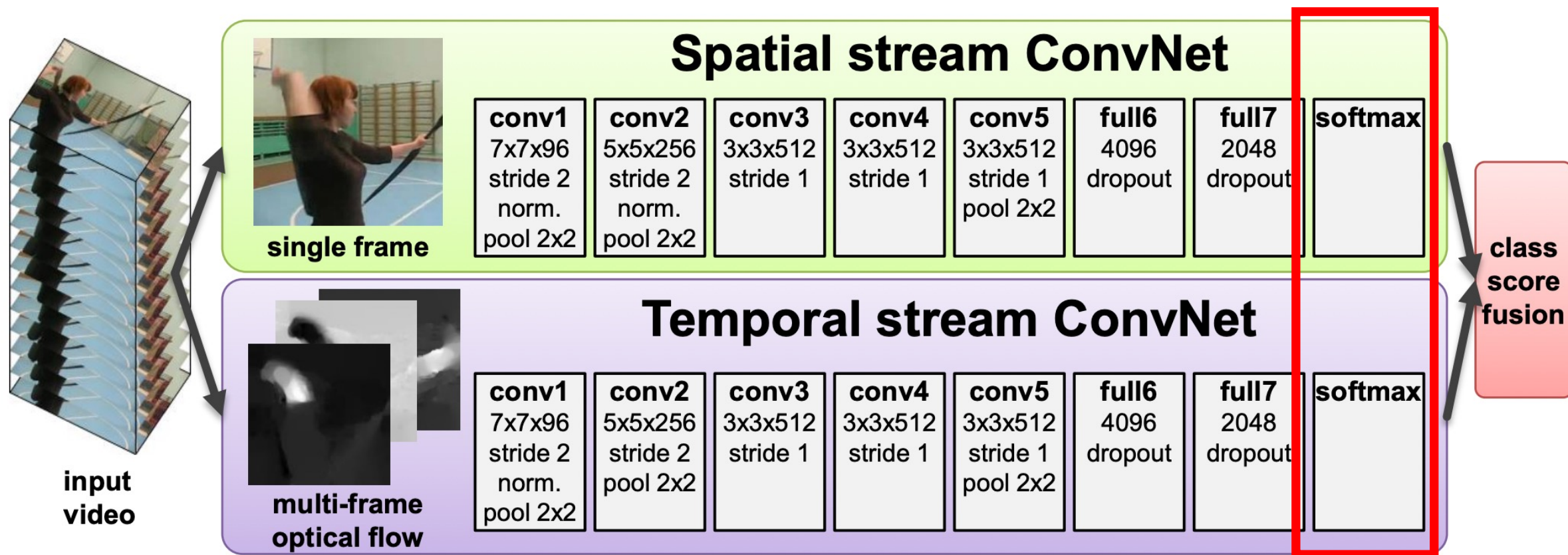


Must be trained on video datasets



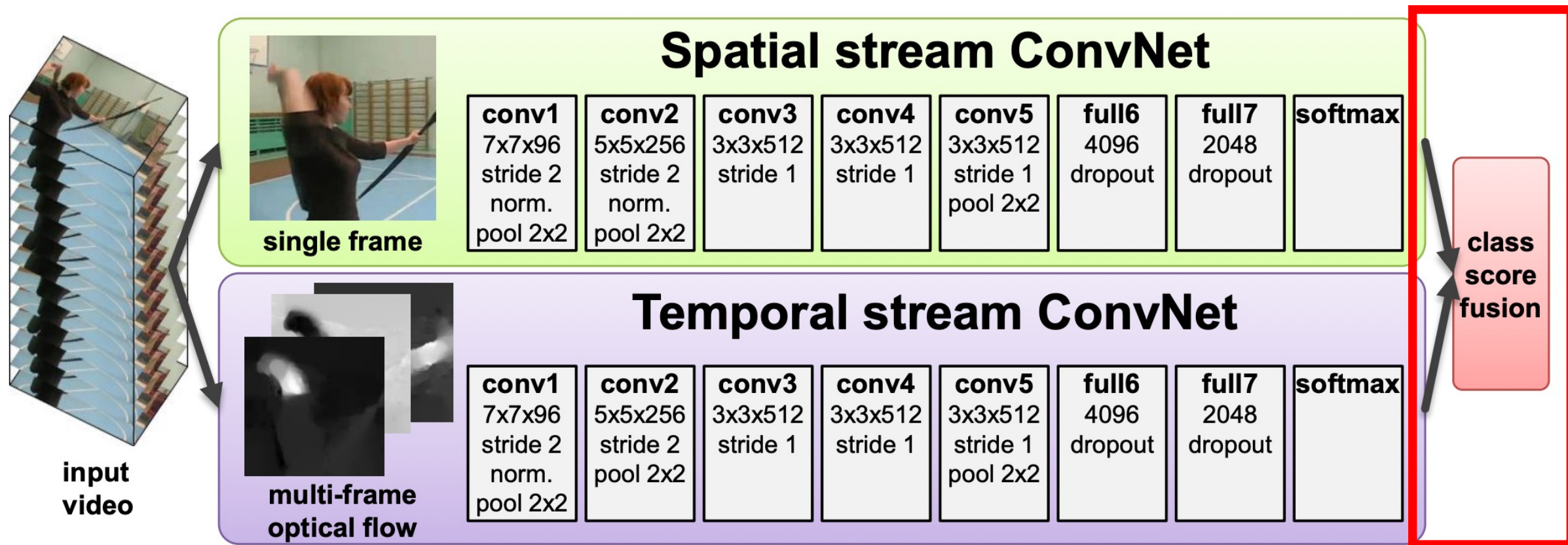
# Two-Stream Architecture

Both ConvNets learn to output a class score



# Two-Stream Architecture

Both class scores can be fused (e.g., averaging or using an SVM)

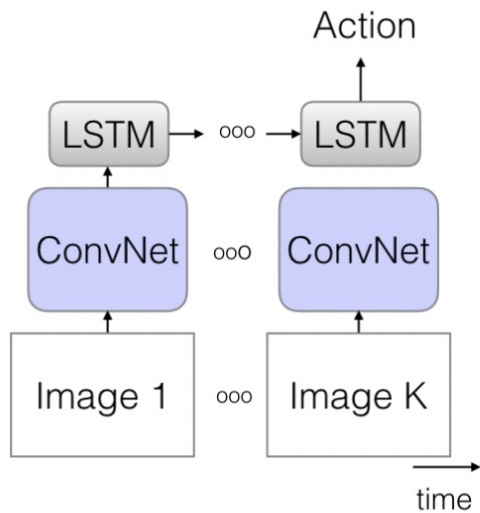


# Two-Stream Architecture: Limitations

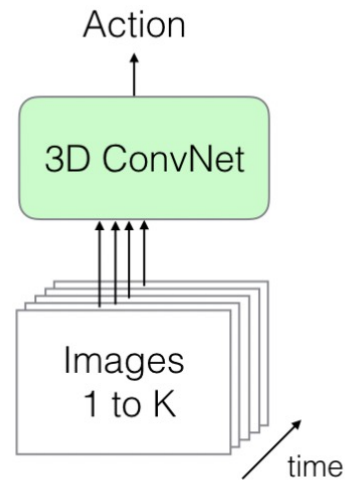
- Successful training requires many videos which the architecture resource-hungry and time-consuming
- Does not capture long-term temporal information (duration determined by sequence duration used for the temporal stream ConvNet)
- Motion representation is limited by the assumptions of optical flow (e.g., constant appearance and smooth flow between frames)

# Approaches to Capture Temporal Information: Can Also Mix Basic Approaches

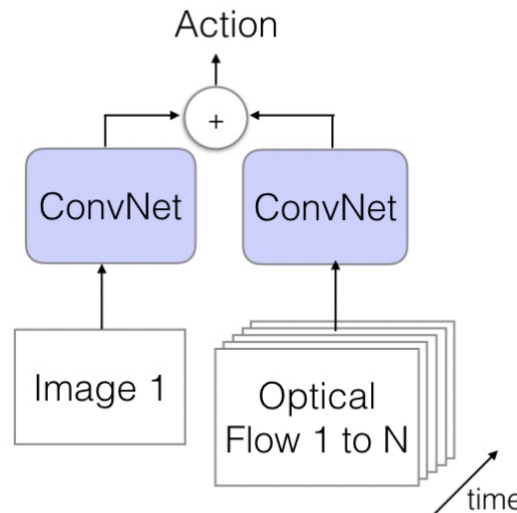
a) LSTM



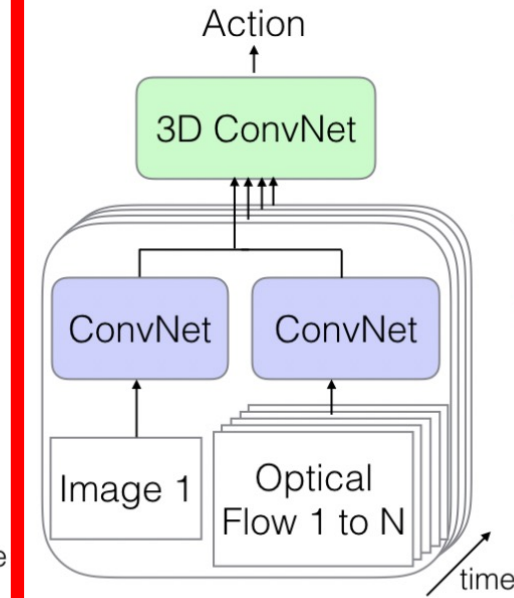
b) 3D-ConvNet



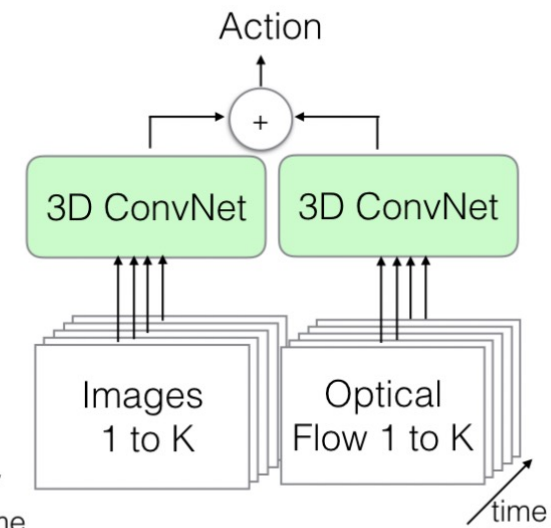
c) Two-Stream



d) 3D-Fused Two-Stream



e) Two-Stream 3D-ConvNet



# Video Classification: Today's Topics

- Problem
- Applications
- Datasets
- Evaluation metric
- Computer vision models

A dark gray background with a central circular glow. The glow is a gradient from light gray in the center to dark gray at the edges. The text "The End" is centered within this glow. The entire scene is framed by a white film strip border with rectangular sprocket holes on the left and right sides.

*The End*