

# Beyond Self-Supervised Representation Learning

Andrew Zisserman

August 2020

# What is Self-Supervision?

- A form of unsupervised learning where the data provides the supervision
- In general, withhold some part of the data, and task the network with predicting it
- This defines a **pre-text task** (or a **proxy loss**), and the network is forced to learn what we really care about, e.g. a semantic representation, in order to solve it
- We can also train networks for tasks directly, beyond learning data representations

# Outline

Self-supervised learning in three parts:

1. Where are we now with representation learning?
2. Beyond representation learning – applicable tasks
3. Roadmap – the three phases of self-supervised learning

# **Part I**

**Where are we now with  
representation learning?**

# “Classical” Self-supervised learning

1. **Image representation**: self-supervised training on **ImageNet** using a proxy task
  2. Supervised training of network for **downstream task** either by linear probe or initialization for fine tuning, e.g. for PASCAL VOC object category detection
- Example **proxy tasks**: Context, Jigsaw, Colourization, Exemplars, RotNet, Clustering, CPC, SimCLR, MoCo, BYOL
  - **Surpass performance of strong supervision (training with class labels)** on a number of downstream tasks, e.g.
    - PASCAL VOC segmentation, object detection
    - NYU depth, ...

# “Classical” Self-supervised learning – video

1. **Video representation**: self-supervised training on **Kinetics** using a proxy task (only visual domain)
  2. Supervised training of network for **downstream task** either by linear probe or fine tuning, e.g. for Action classification on UCF-101 or HMDB51
- Example **proxy tasks**: Slowness, Shuffle&Learn, Order, Odd-One-Out, AoT, ST-Puzzle, DynamoNet, DPC, CBT, SpeedNet, MemDPC, CoCLR
  - **Approaching the performance of strong supervision** on downstream tasks

## **Part II**

# **Beyond representation learning – applicable tasks**

# Outline

Traditional: learn data representation with proxy task using self-supervision, then linear probe or finetune for downstream task using supervision

Instead, train for an **applicable task** directly using self-supervision

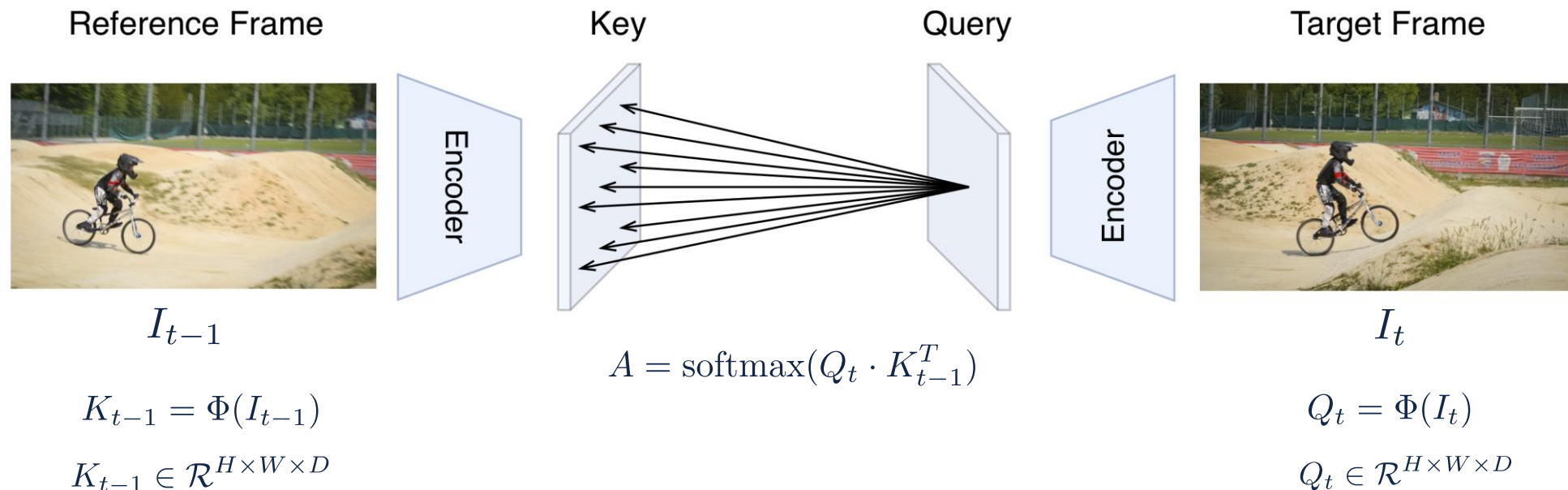
Illustrate with three example tasks on video:

1. Object tracking in videos
2. Audio-visual joint embedding and localization
3. Obtain discrete audio-visual objects

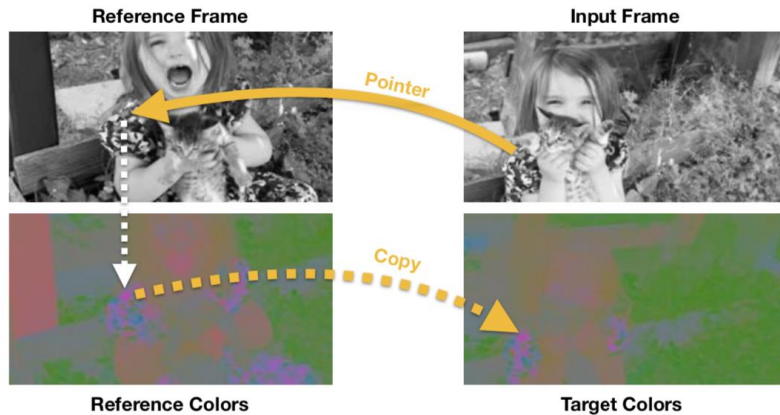


# Applicable task 1: Self-supervised Learning for Video Object Tracking

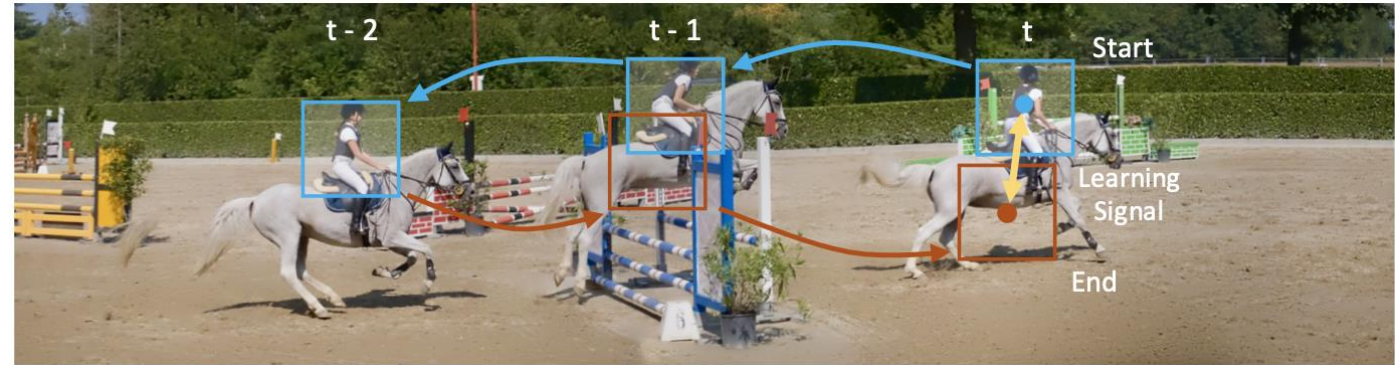
- Tracking can be solved by learning the pixelwise correspondence between consecutive frames
- Use an attention mechanism between spatial features of each frame to determine a soft correspondence



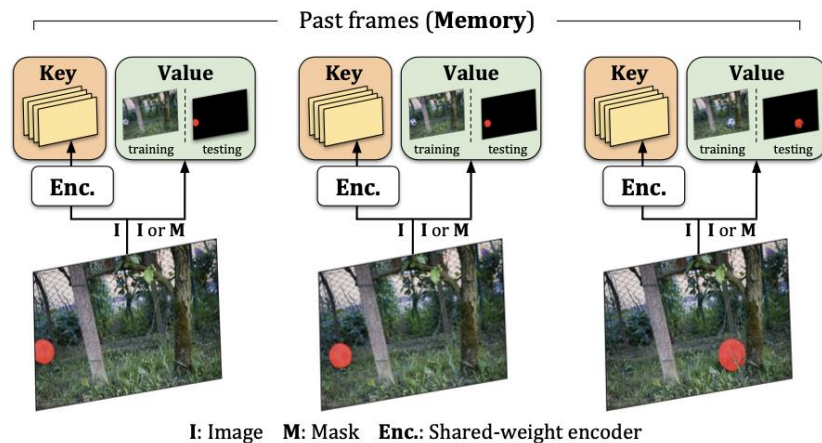
# Applicable task 1: Self-supervised Learning for Video Object Tracking



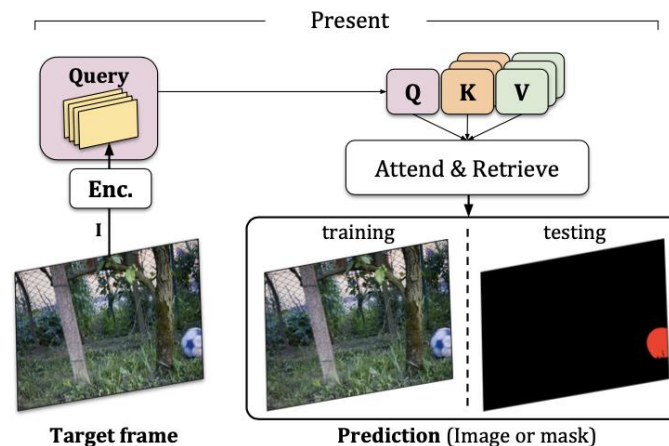
Tracking Emerges by Colorizing Videos  
[Vondrick *et al.* ECCV 2018]



Learning Correspondence from the Cycle-consistency of Time  
[Wang, Jabri & Efros, CVPR2019]



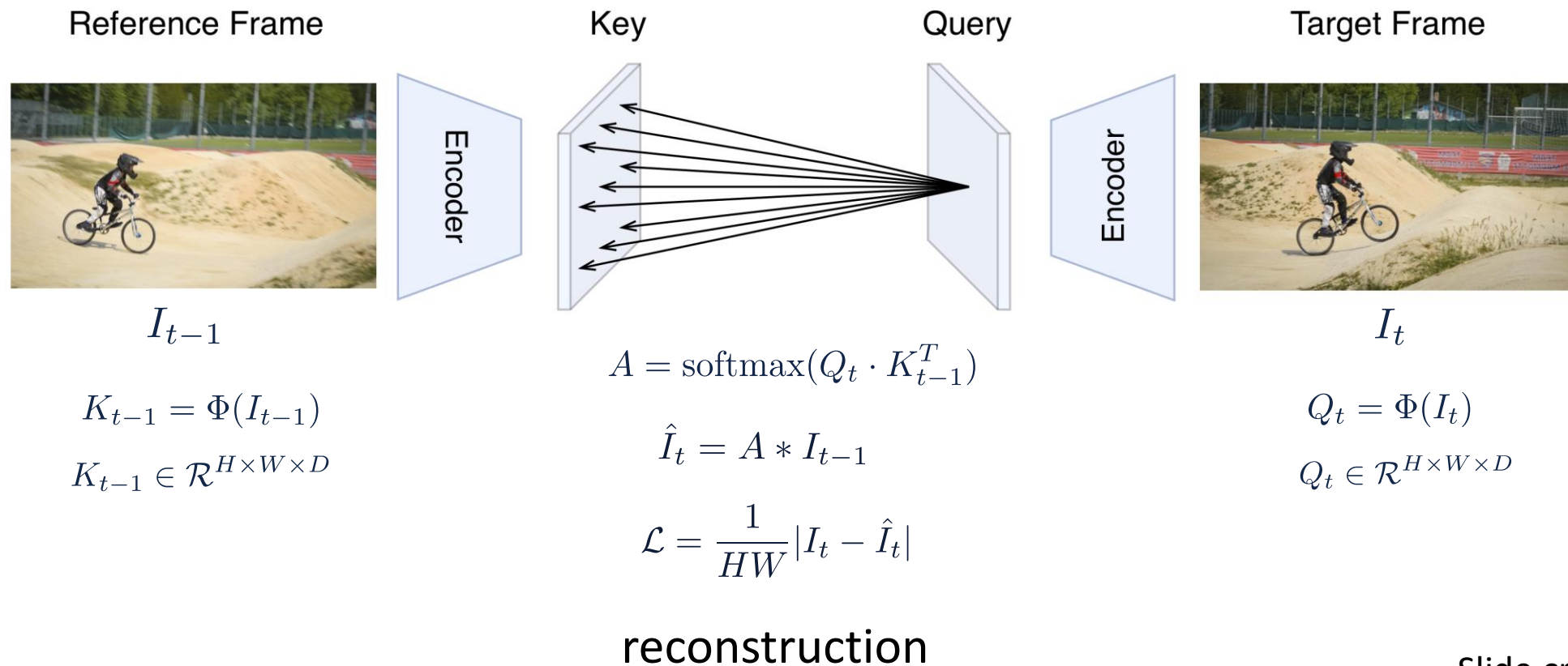
I: Image M: Mask Enc.: Shared-weight encoder



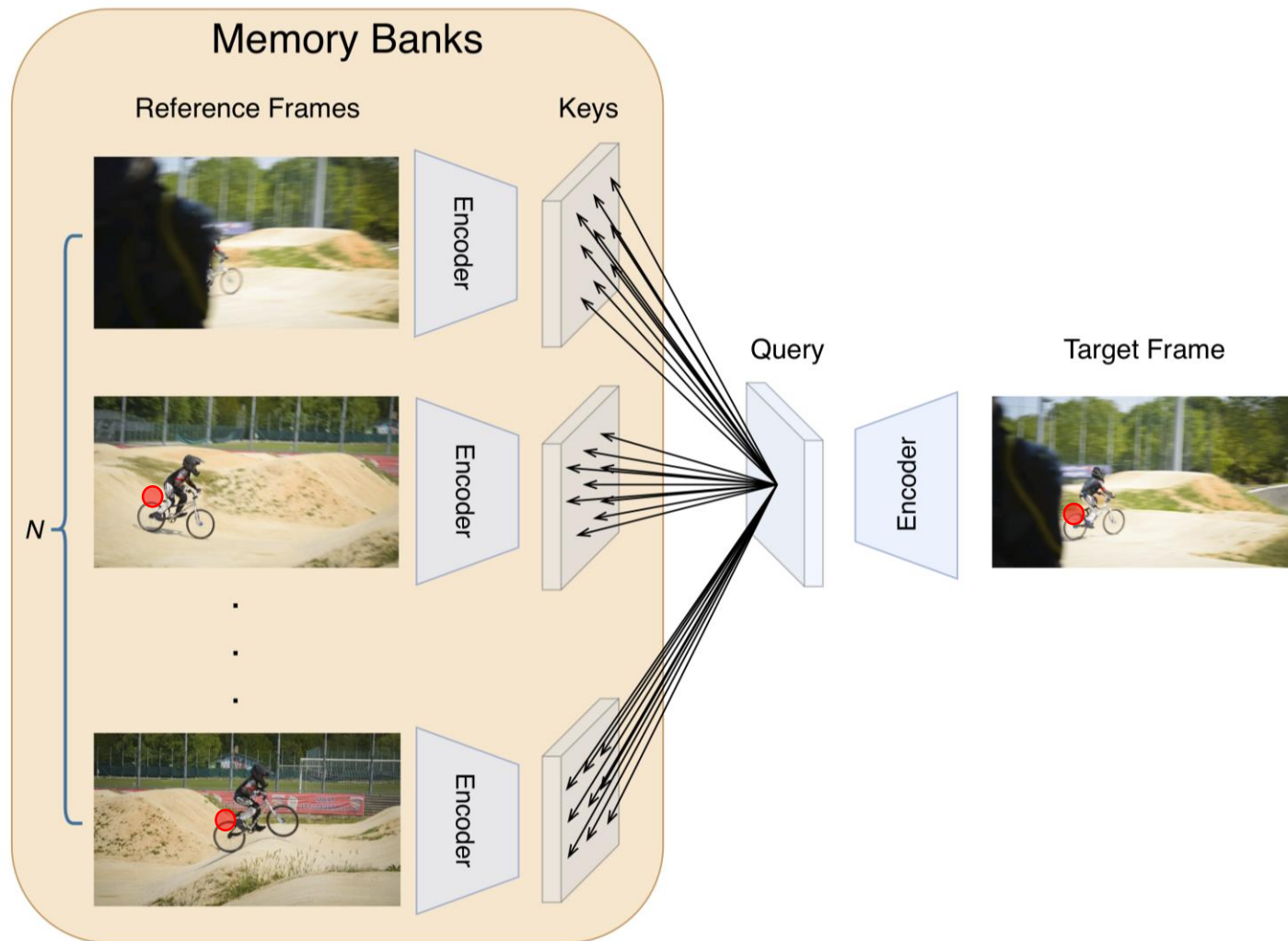
MAST: A Memory-Augmented  
Self-supervised Tracker  
[Lai, Lu & Xie, CVPR2020]

# Applicable task 1: Self-supervised Learning for Video Object Tracking

- Use an attention mechanism between spatial features of each frame to determine a soft correspondence
- Learn by reconstructing a target frame by copying pixels from a previous frame or by cycle consistency



# Memory-augmented Self-supervised Tracking



- Construct the memory bank with multiple reference frames, affinity matrix becomes:

$$A \in \mathcal{R}^{HW \times HW N}$$

$$\hat{I}_t = A * [I_{t-N}, \dots, I_{t-1}]$$

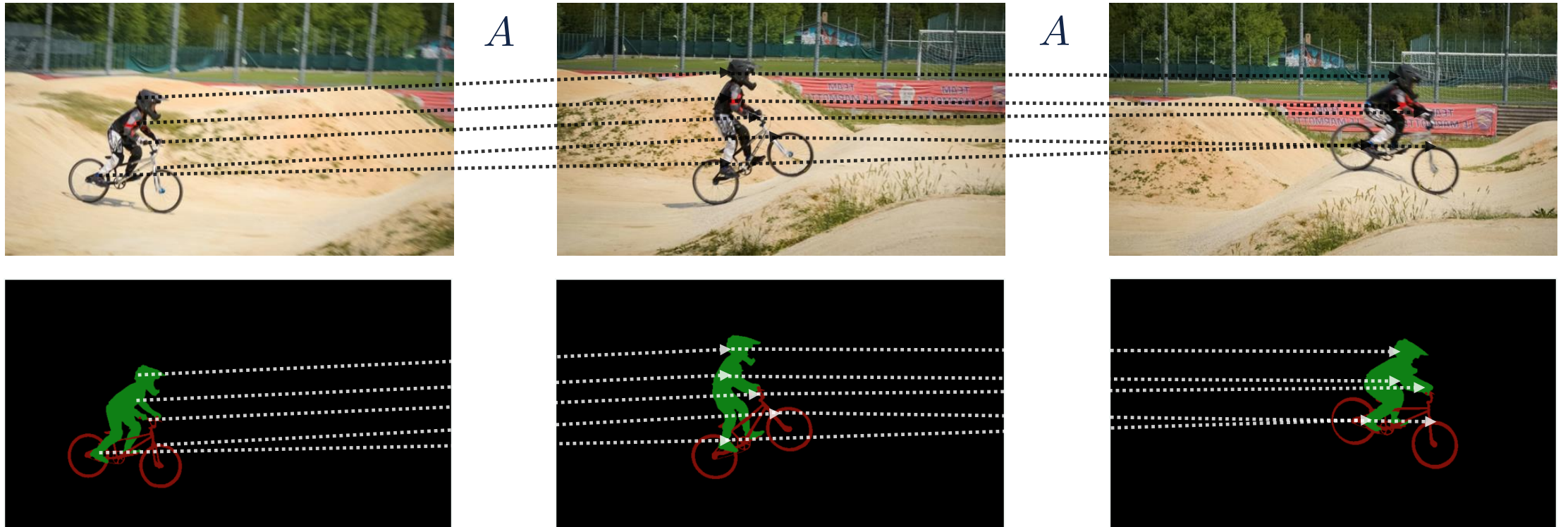
- Effectively handle the occlusion problems, reducing the tracker drift.



# How to achieve Self-supervised Tracking ?

- Propagate instance masks from previous frames:

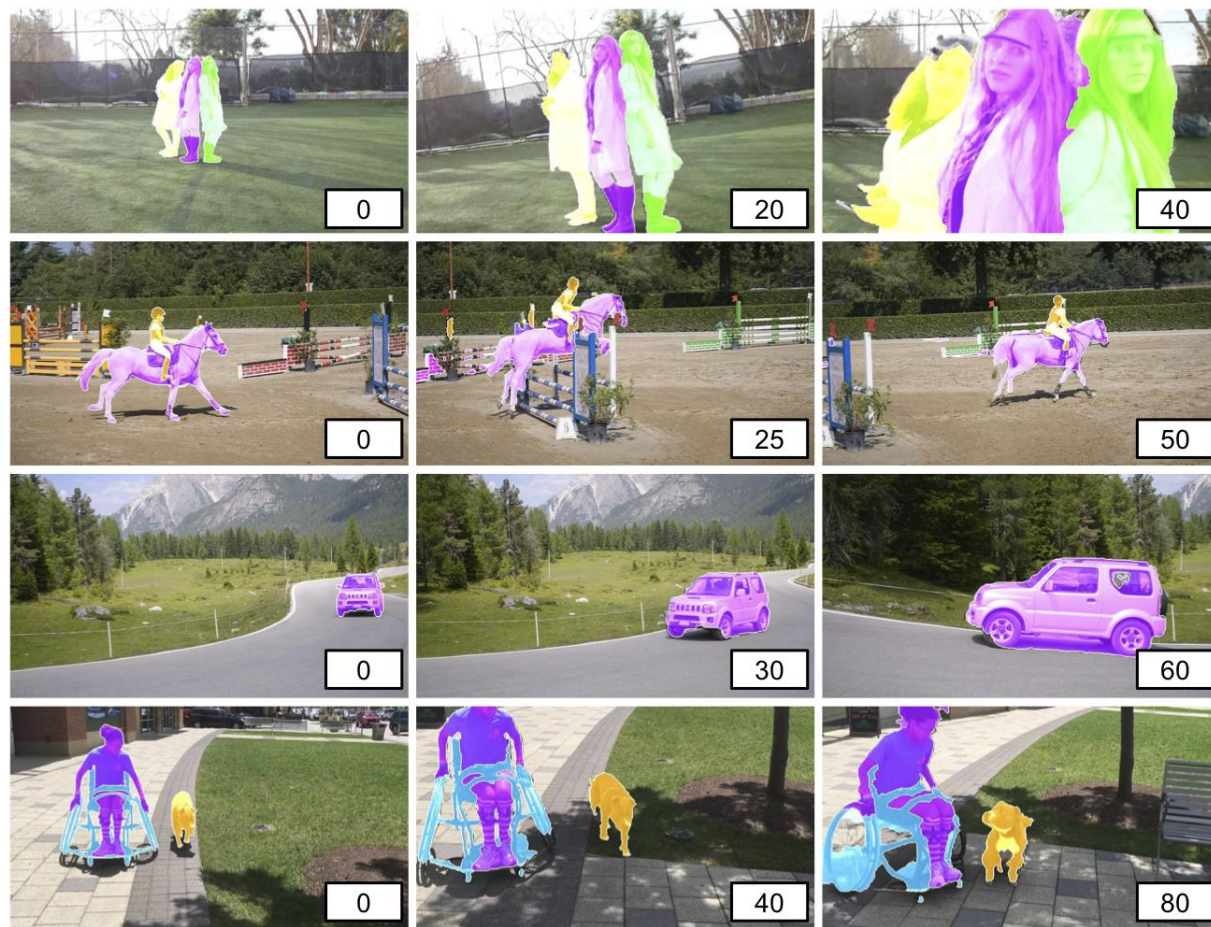
$$\hat{M}_t^i = \sum_j A_t^{ij} M_{t-1}^j$$



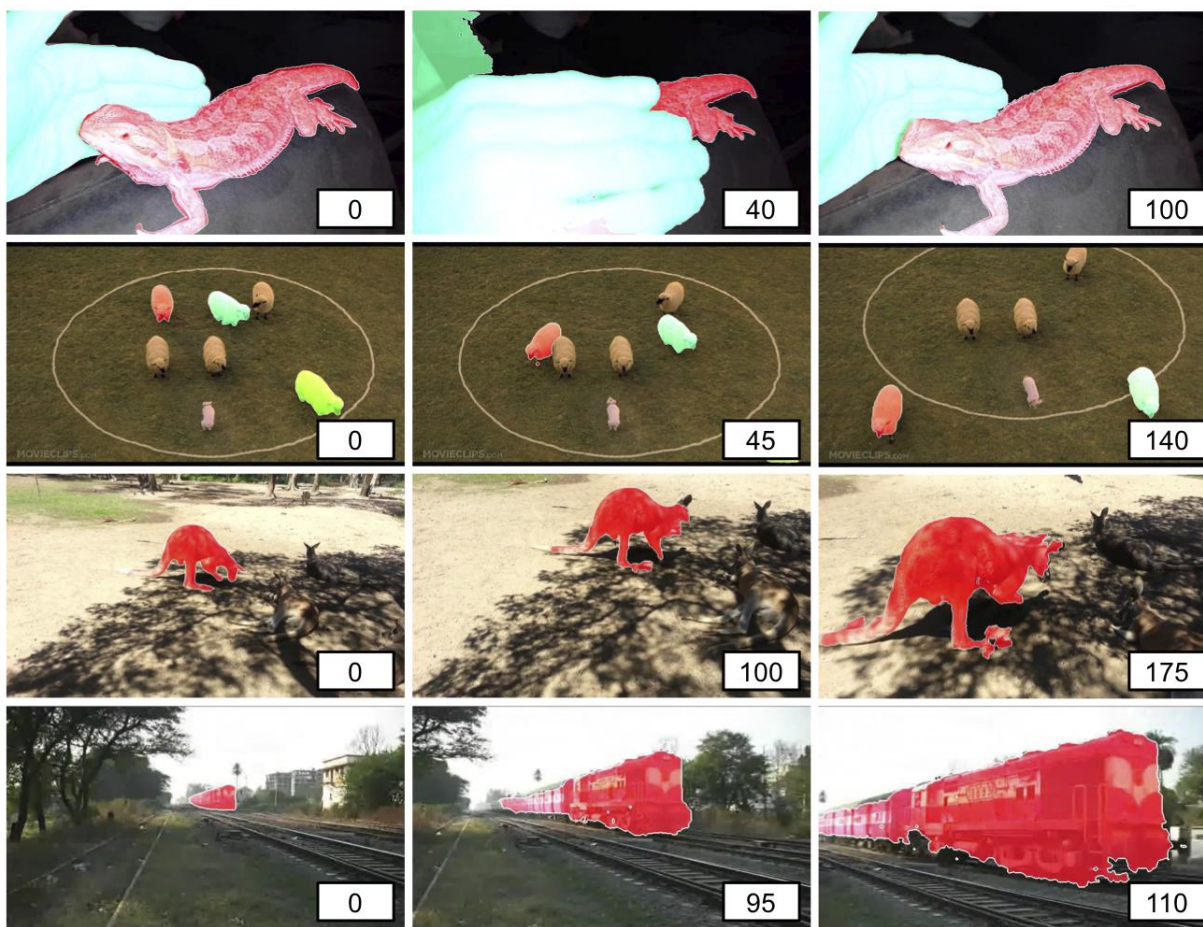


# Qualitative Results

DAVIS-2017

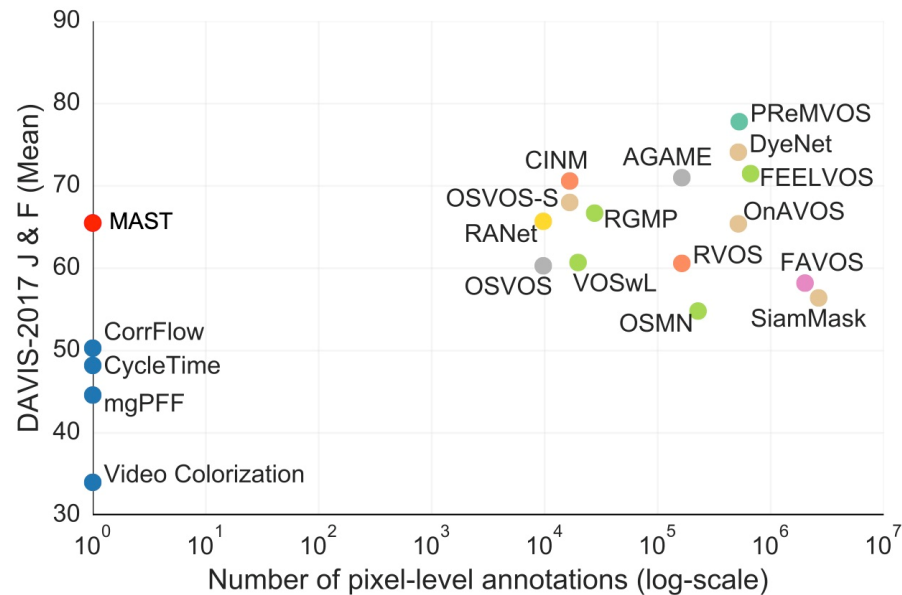


YouTube-VOS



# What has been achieved ?

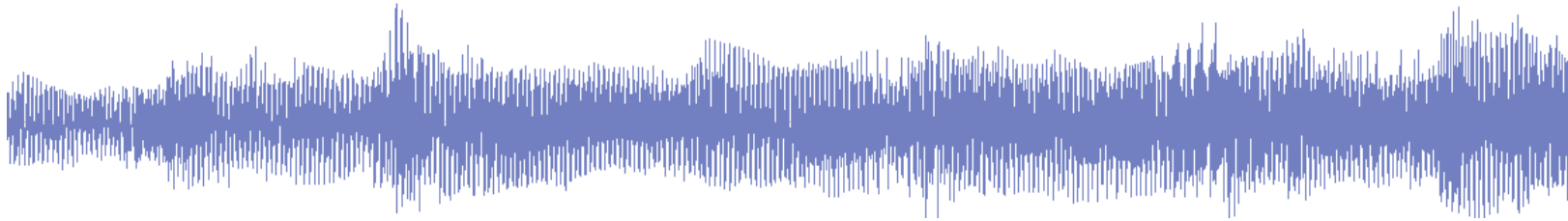
- Benchmark on the public DAVIS Video Segmentation Dataset.
- Over the last two years, self-supervised approaches have shown great promise on the task of dense tracking, outperforming many supervised ones, trained with millions of expensive pixel-wise segmentation annotations.





# Audio-Visual Co-supervision

**Objective:** use vision and sound to learn from each other

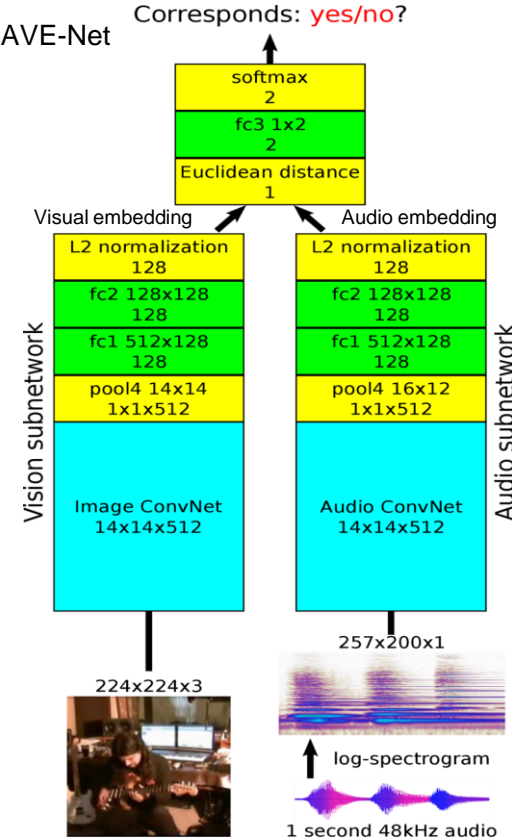


- Sound and frames are (i); synchronized, and (ii) semantically consistent
- Two types of proxy task:
  1. Predict audio-visual **correspondence**
  2. Predict audio-visual **synchronization**

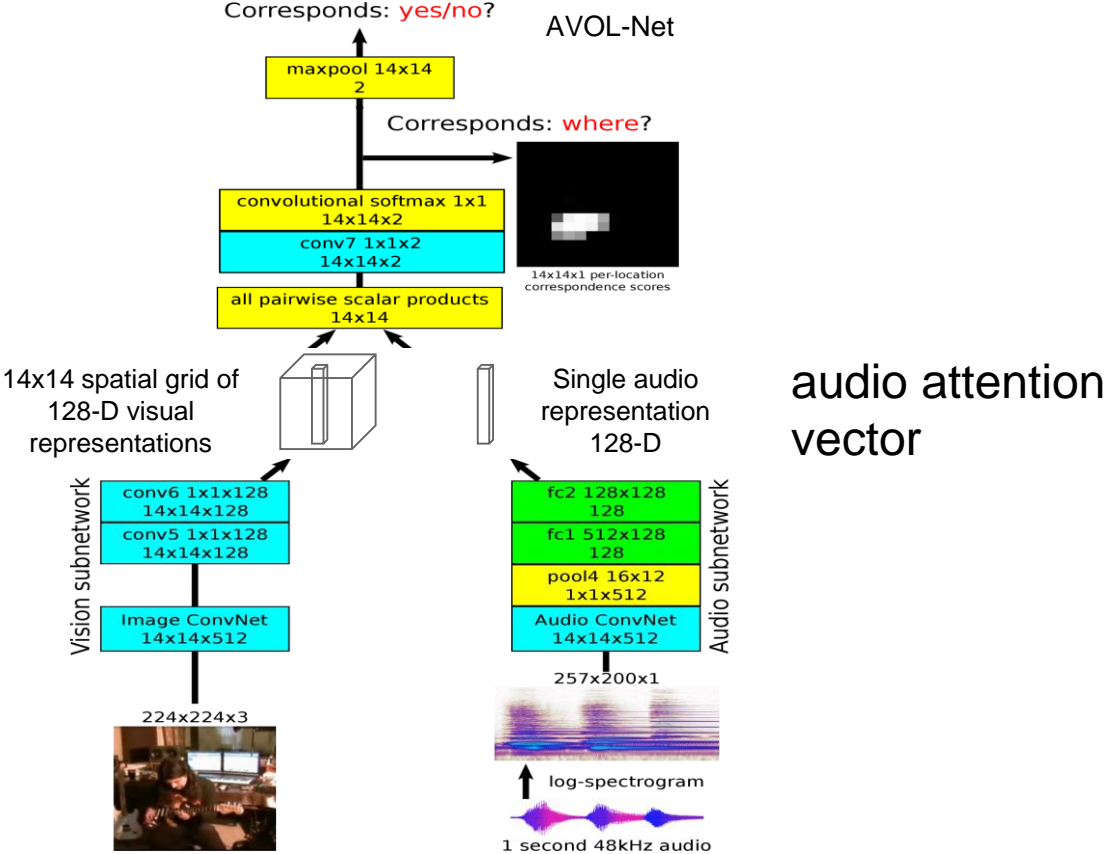


# Applicable task 2: Audio-visual joint embedding and localization

## Joint embedding



## Joint embedding and localization



# Audio-Visual Co-supervision

Train a network to predict if **image** and audio clip correspond

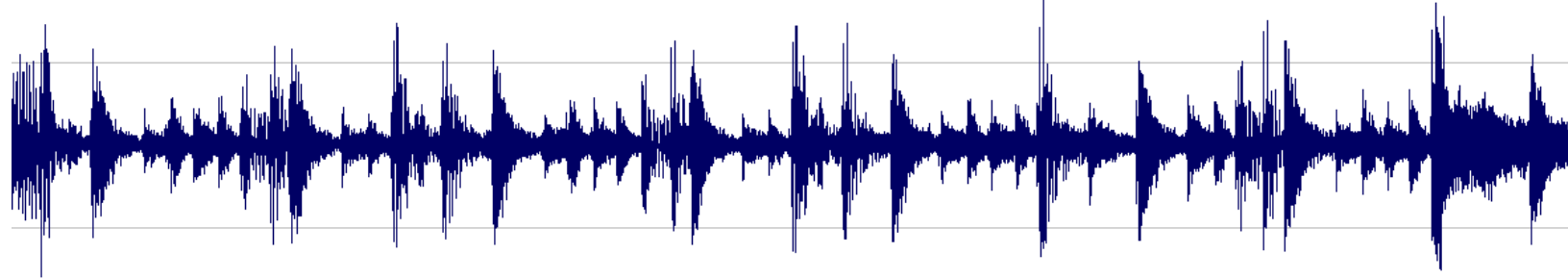


Correspond?

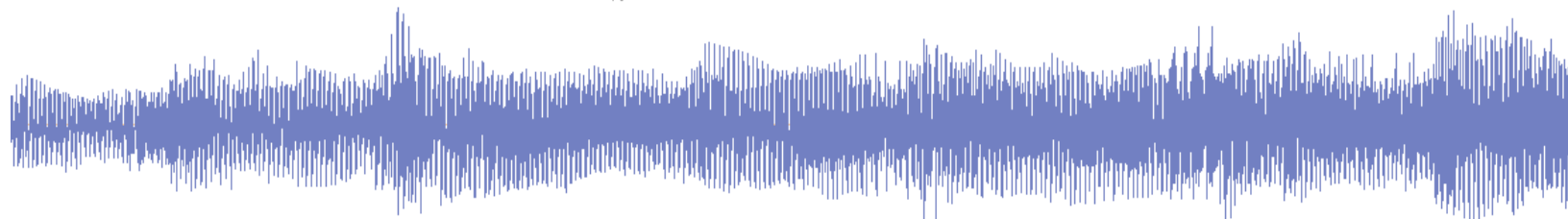


# Audio-Visual Correspondence

drum



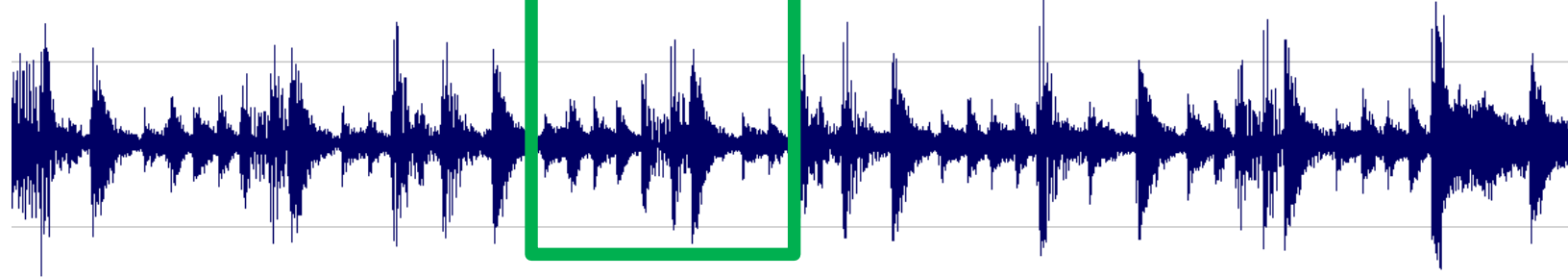
guitar



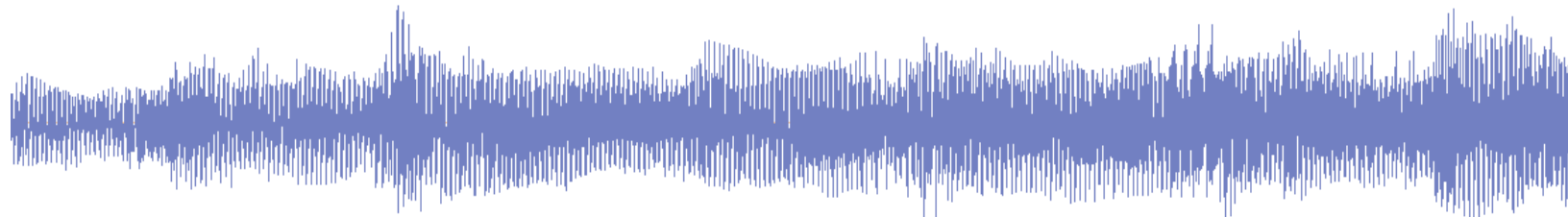
# Audio-Visual Correspondence

positive

drum



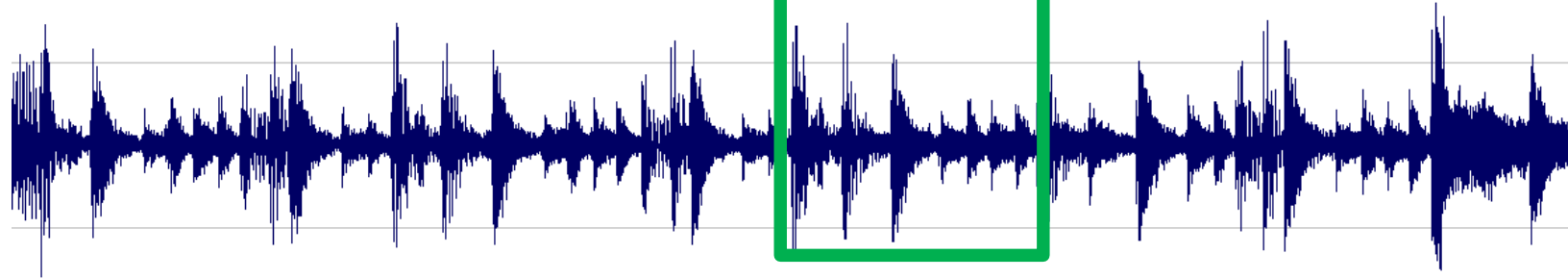
guitar



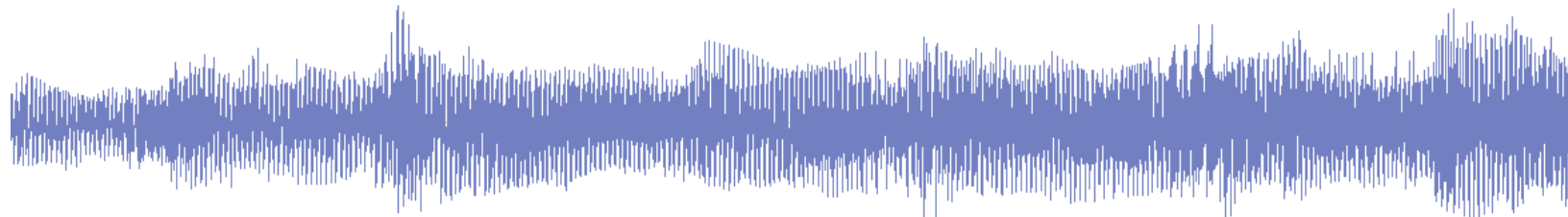
# Audio-Visual Correspondence

positive

drum

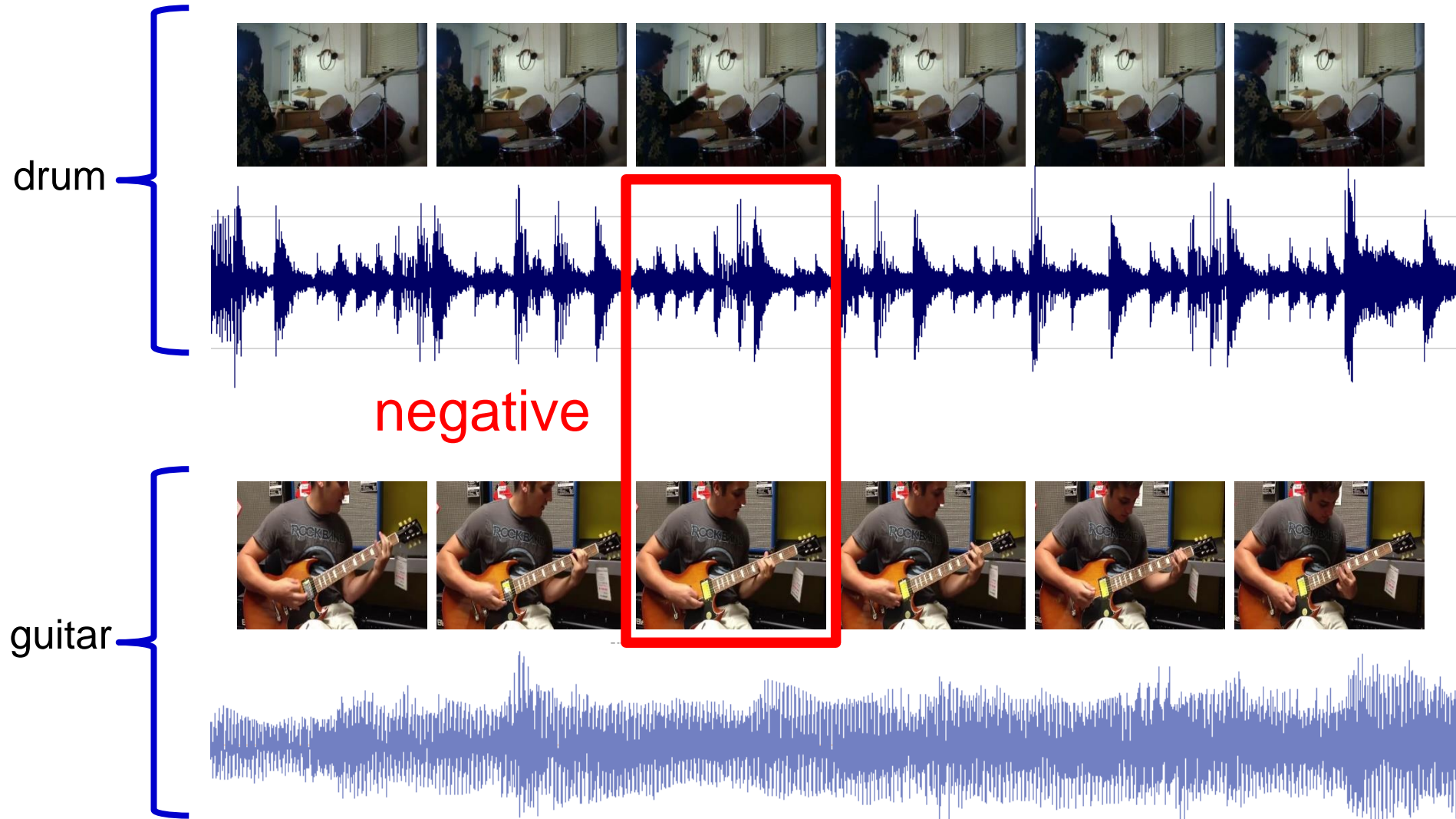


guitar

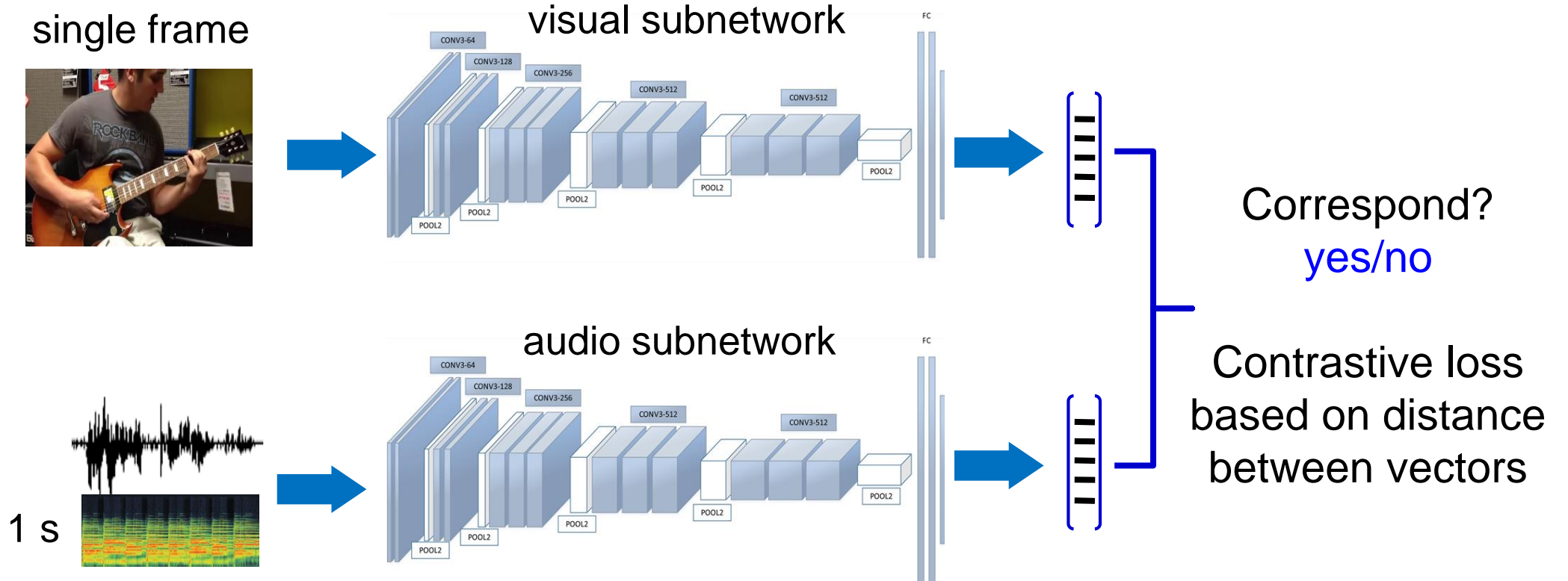




# Audio-Visual Correspondence



# Audio-Visual Embedding (AVE-Net)



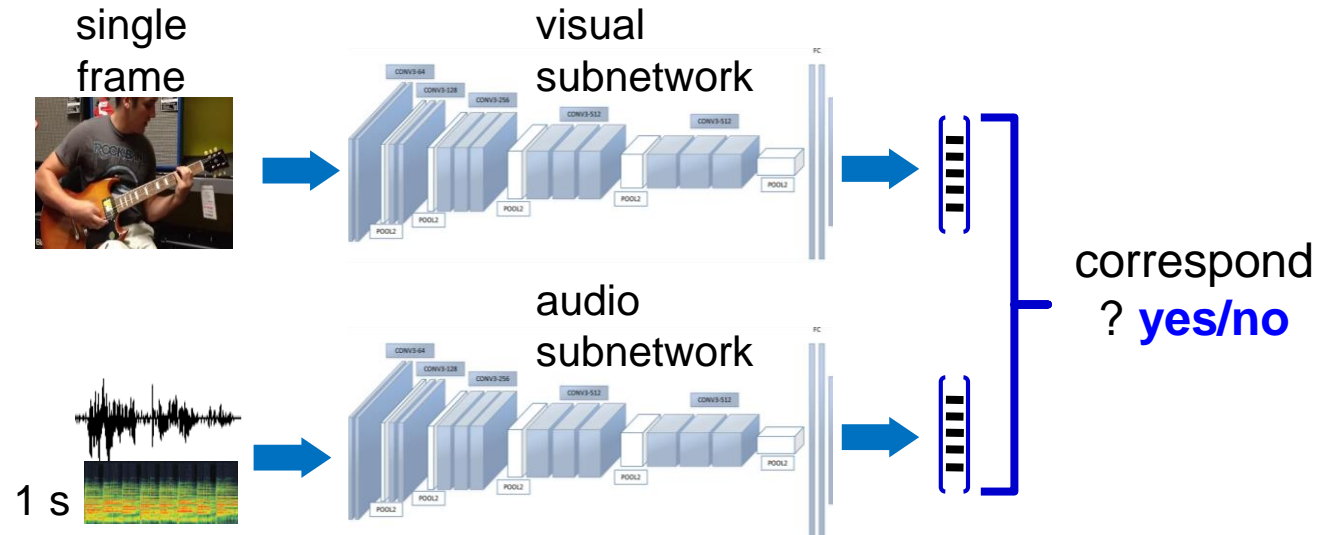
**Distance between audio and visual vectors:**

- **Small:** AV from the same place in a video (**Positives**)
- **Large:** AV from different videos (**Negatives**)

Train network from scratch

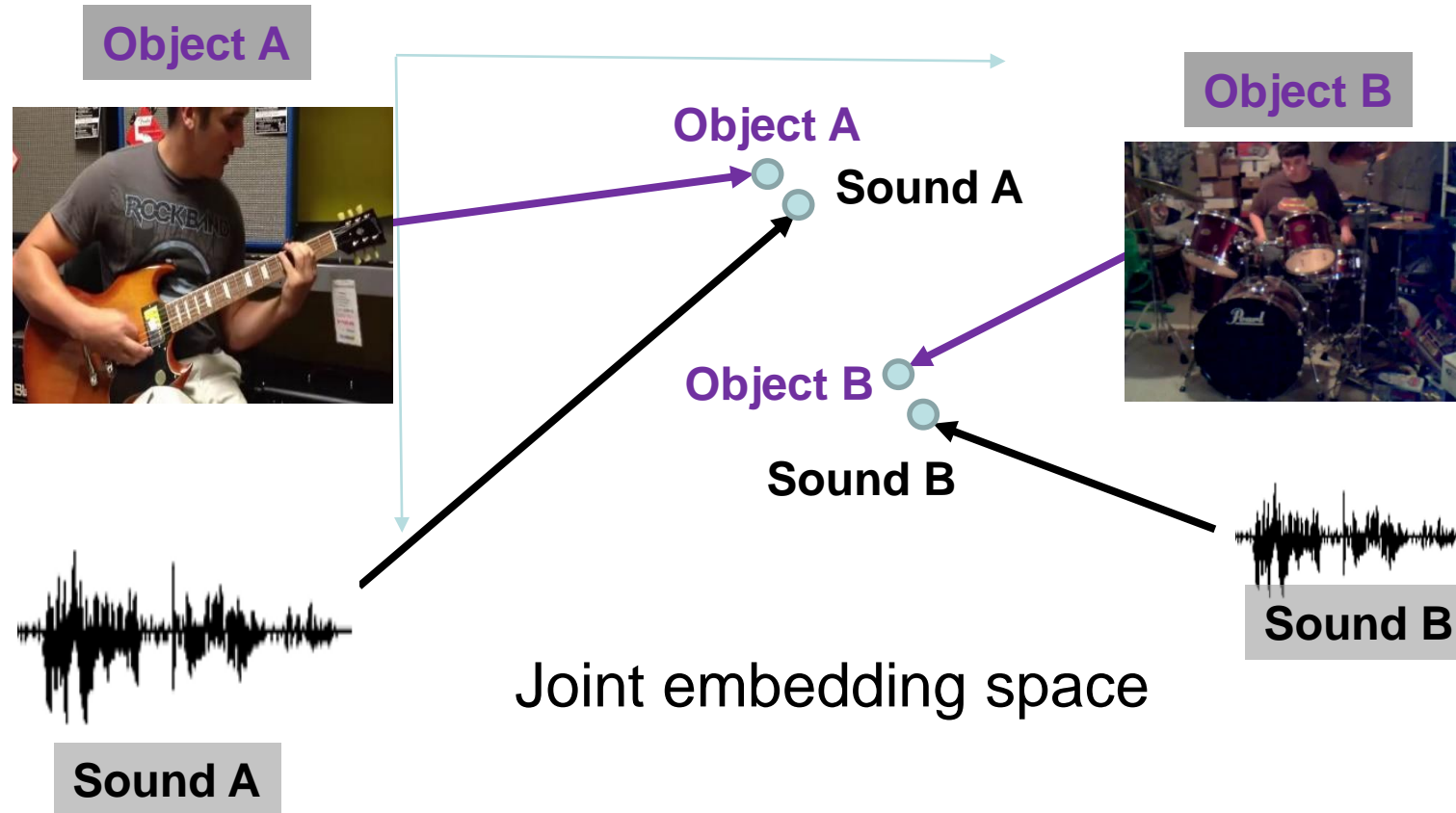
# What has been learnt?

- Good representations
  - Visual features
  - Audio features
- Intra- and cross-modal retrieval
  - Aligned audio and visual embeddings





# Joint Embedding



# Query on audio, retrieve image

- Audio to Vision

Query



Audio-to-image retrieval: Top 5 retrieved results



# Query on image, retrieve audio

Search in 200k video clips of AudioSet

Query  
frame

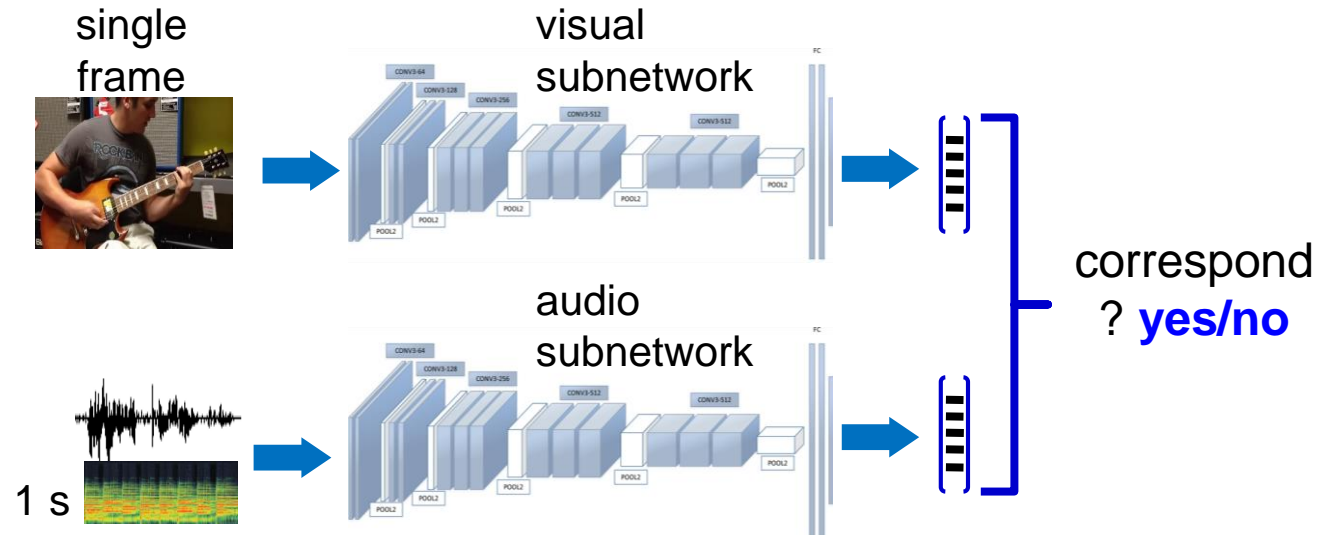


Top 10 ranked audio clips



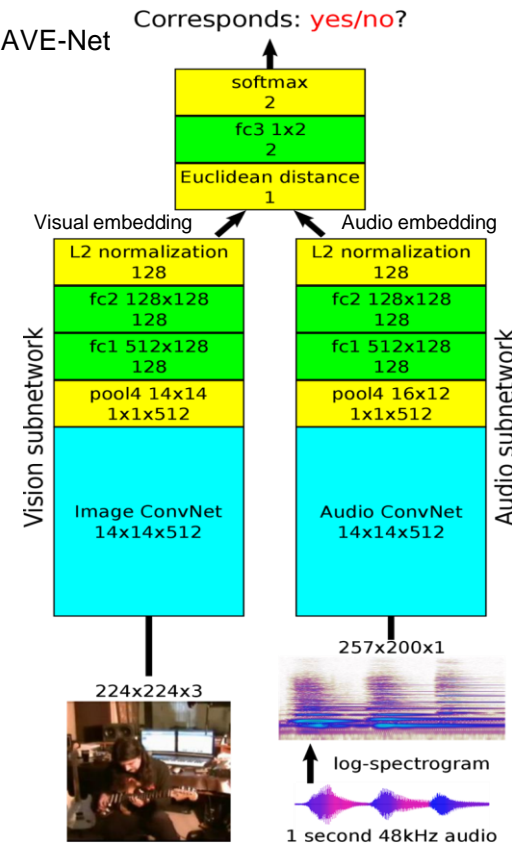
# Audio-visual joint embedding and localization

- Good representations
  - Visual features
  - Audio features
- Intra- and cross-modal retrieval
  - Aligned audio and visual embeddings
- “What is making the sound?”
  - Learn to localize objects that sound

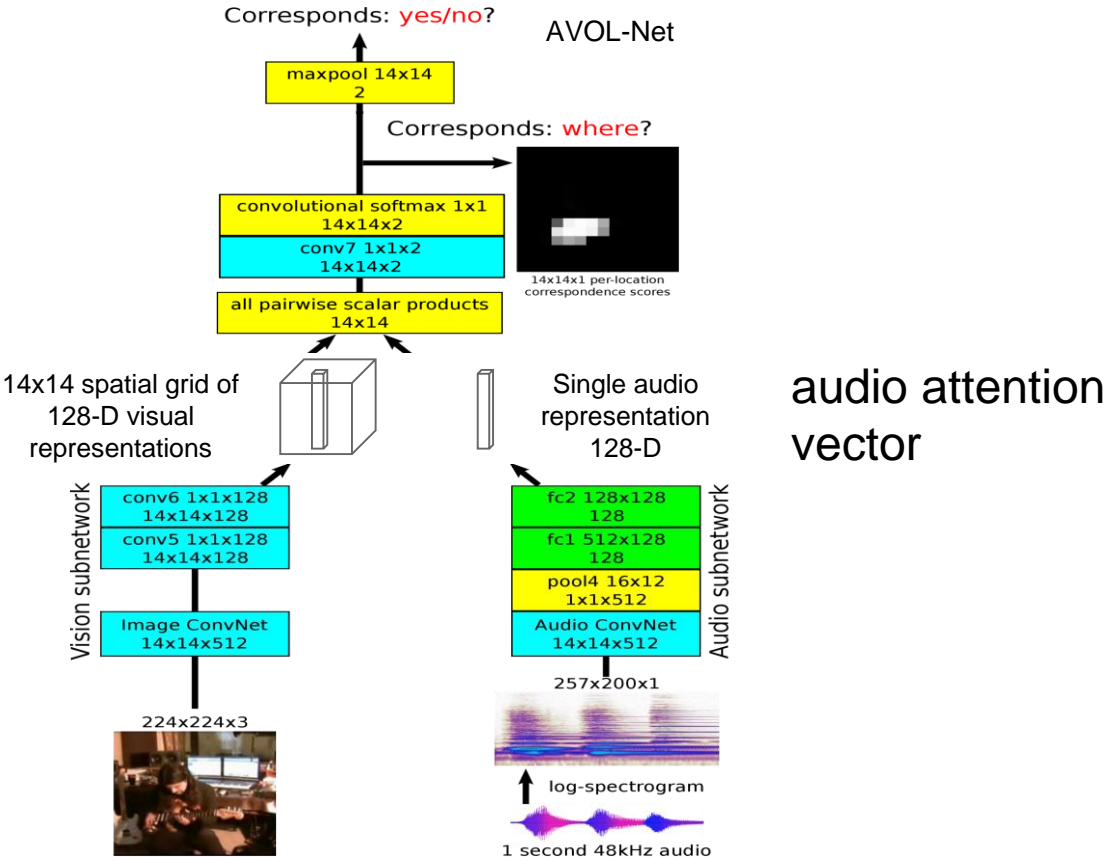


# Applicable task 2: Audio-visual joint embedding and localization

## Joint embedding



## Joint embedding and localization



# Objects that Sound: object localization

Input: audio and video frame



frame

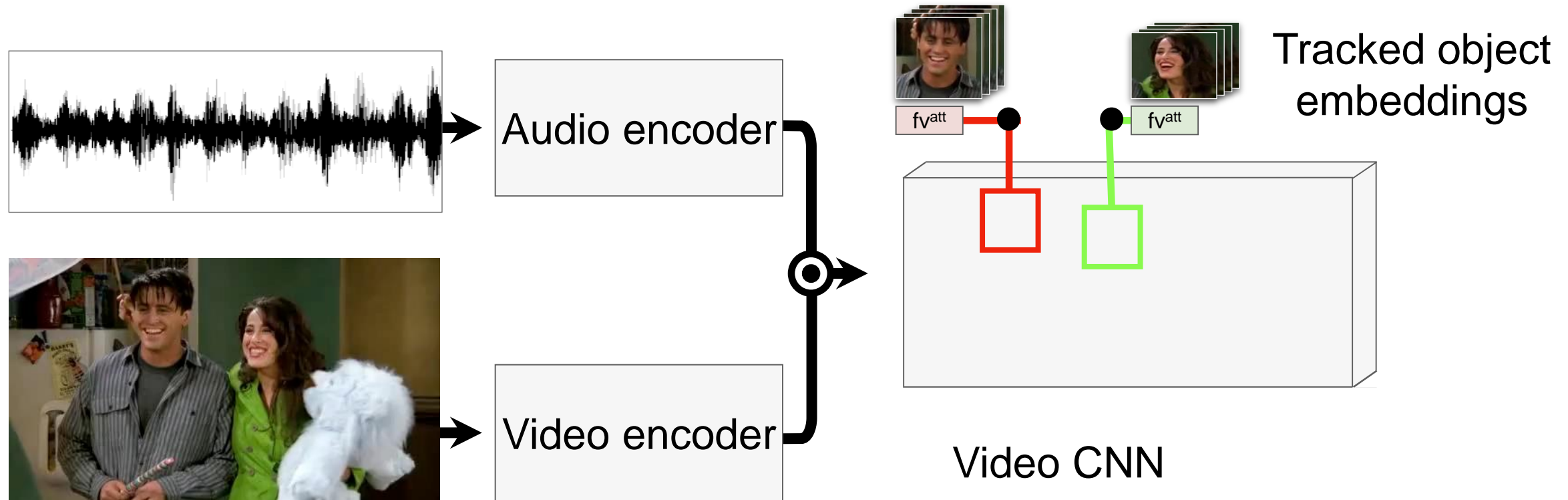
frame+heatmap

heatmap

- Frame by frame
- No motion information
- No memory
- No smoothing

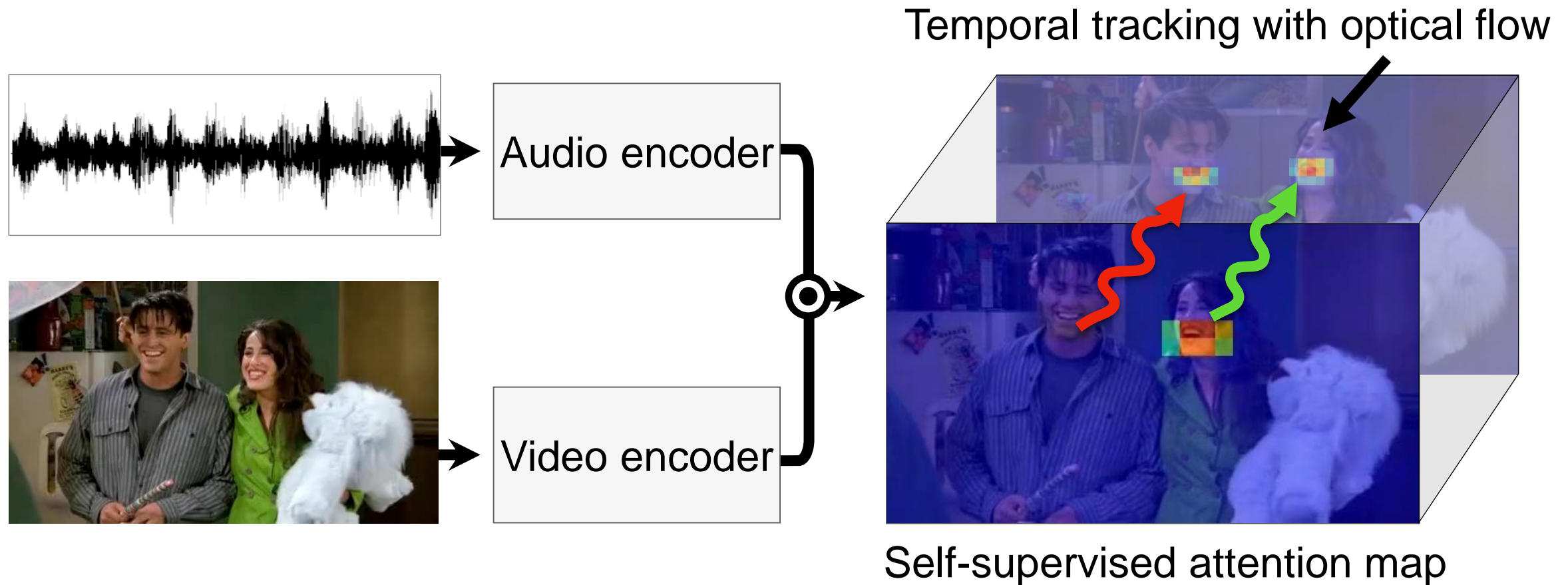


# Applicable task 3: : obtain discrete audio-visual objects



Self-Supervised Learning of Audio-Visual Objects from Video  
T. Afouras, A. Owens, J. S. Chung, A. Zisserman, ECCV 2020

# Applicable task 3: : obtain discrete audio-visual objects



Self-Supervised Learning of Audio-Visual Objects from Video  
T. Afouras, A. Owens, J. S. Chung, A. Zisserman, ECCV 2020



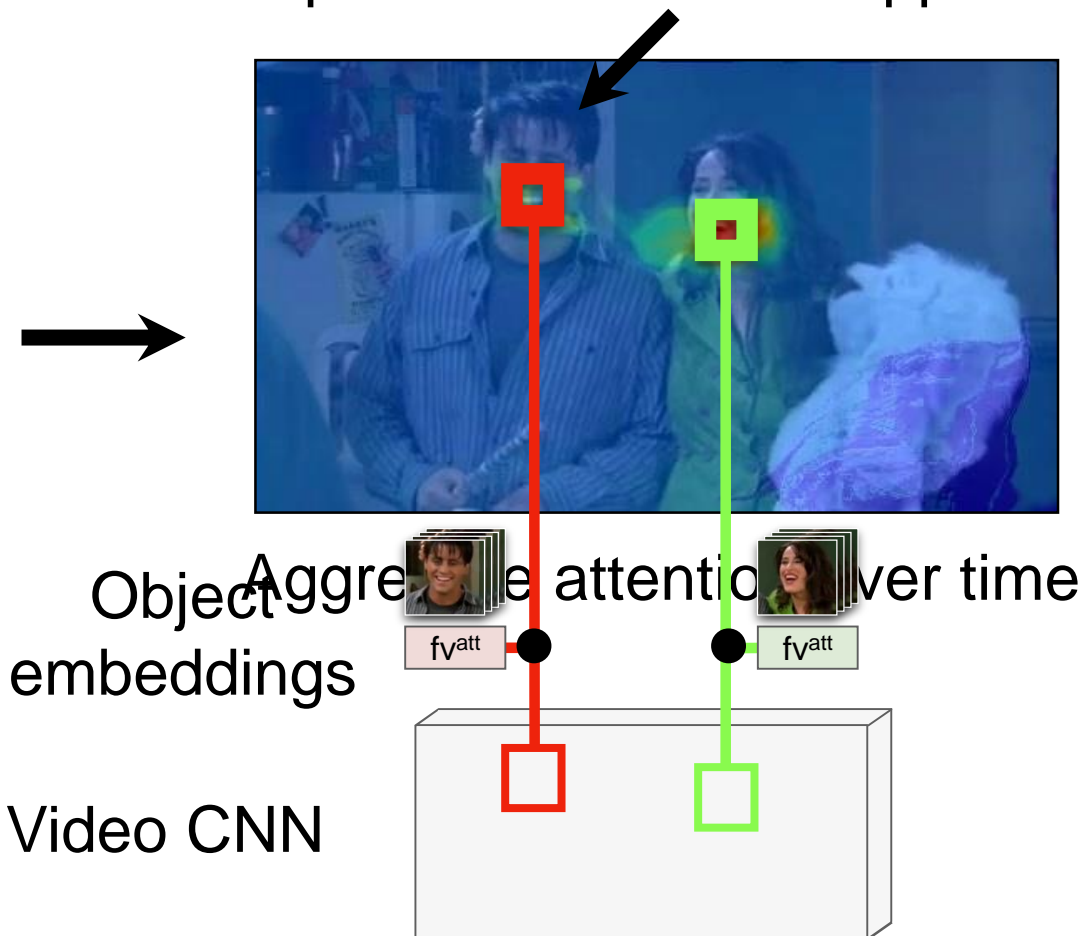
# Applicable task 3: : obtain discrete audio-visual objects

Temporal tracking with optical flow



Self-supervised attention map

Find peaks + non-max suppression



Self-Supervised Learning of Audio-Visual Objects from Video  
T. Afouras, A. Owens, J. S. Chung, A. Zisserman, ECCV 2020

# Learning the attention maps

- Contrastive loss:
- **Positive samples:** in sync
- **Negative samples:** out of sync (with temporal offset)

positive



See also: [Chung & Zisserman 2016], [Owens & Efros 2018], [Arandjelović & Zisserman 2018], [Korbar et al. 2018]

# Learning the attention maps

- Contrastive loss:
- **Positive samples:** in sync
- **Negative samples:** out of sync (with temporal offset)

negative



See also: [Chung & Zisserman 2016], [Owens & Efros 2018], [Arandjelović & Zisserman 2018], [Korbar et al. 2018]

# Learning the attention maps

- Contrastive loss:
- **Positive samples:** in sync
- **Negative samples:** out of sync (with temporal offset)

positive



See also: [Chung & Zisserman 2016], [Owens & Efros 2018], [Arandjelović & Zisserman 2018], [Korbar et al. 2018]

# Learning the attention maps

- Contrastive loss:
- **Positive samples:** in sync
- **Negative samples:** out of sync (with temporal offset)

negative



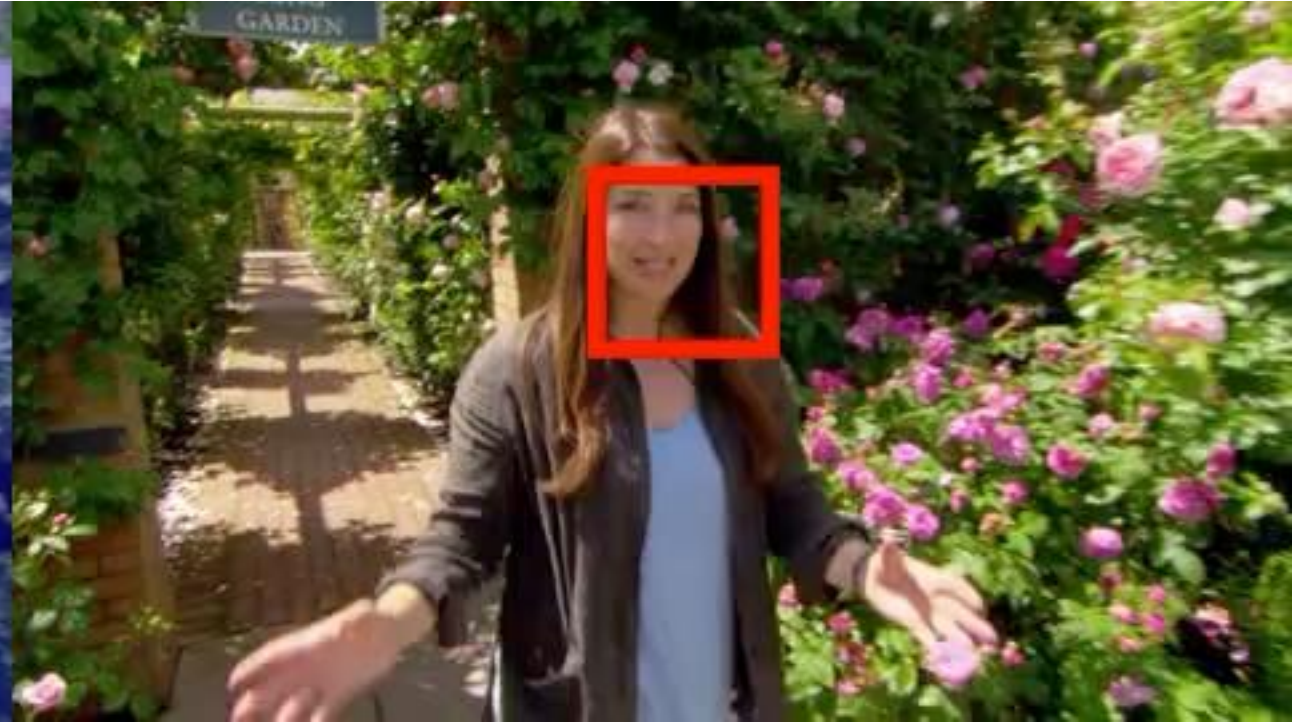


# Audio-Visual Objects: tracking

Examples from the LRS2 dataset



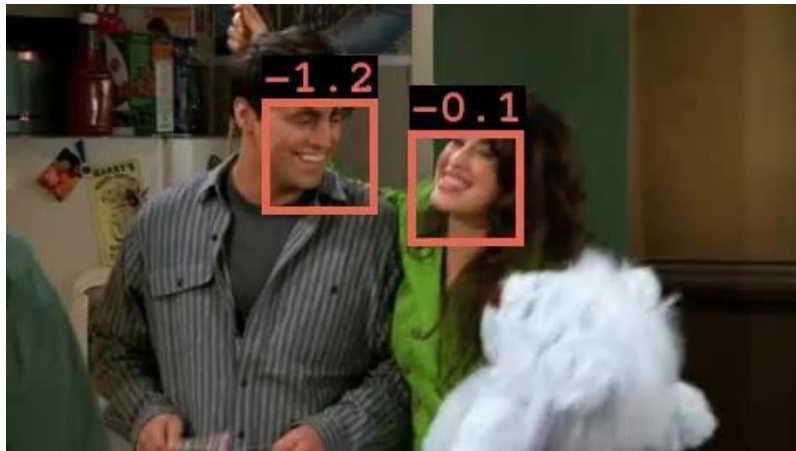
$S_{AV}$  attention map



Audio-visual object

And have tracked visual embeddings for individual objects

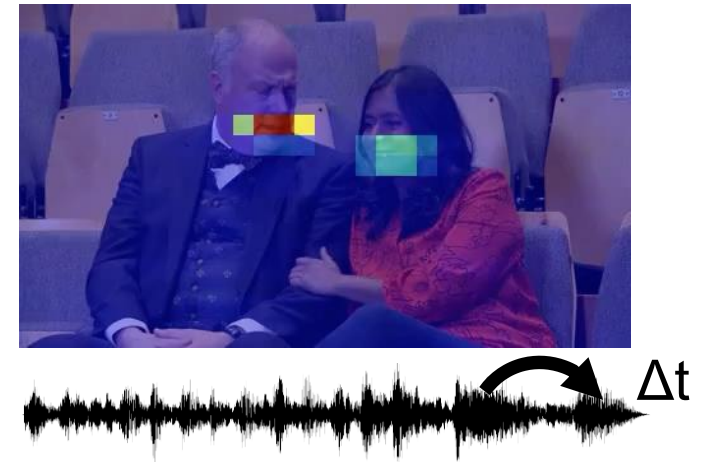
# Applications of audio-visual objects



Active speaker  
detection



Multi-speaker  
source separation



Correcting temporal  
misalignment



# Active Speaker Detection

Examples from the *Friends* series



Blue = active speaker  
Red = inactive speaker



# Adapting to new domains ...

- Since everything is self-supervised, just fine tune



Sesame Street



The Simpsons

# Active Speaker Detection

Examples from *Sesame Street*



Blue = active speaker  
Red = inactive speaker

# Active Speaker Detection

Examples from *The Simpsons*



Blue = active speaker  
Red = inactive speaker

# Summary

- Self-supervision directly for applicable tasks (here discrete audio-visual object extraction)
- Many benefits accrue without having to train for them
  - Visual embedding vector for each object
  - Attention localization from audio
  - Use embedding vector for (more) downstream tasks, e.g. source separation
  - Plug and play for new videos of talking humans
  - Fine tune for non-human (same architecture, same self-supervised proxy)
- Compare to what we don't have to do
  - No two stage: representation learning then downstream
  - No face/head detector required
  - No prior grouping of faces into tracks
  - Video volume processed as a whole, rather than processing each face track

# **Part III**

## **Roadmap: the three phases**



# Phase 1: the “classic” phase

- Replace strong supervision with self-supervision for representation learning
- Goals:
  - develop proxy loss for training an image representation network on ImageNet, evaluate on downstream image tasks
  - develop proxy loss for training a video representation network on Kinetics, evaluate on downstream video tasks
- Drop in replacement for supervised training
- Example proxy tasks for images: Context, Jigsaw, Colourization, Exemplars, RotNet, CPC, SimCLR, MoCo, BYOL
- Example proxy tasks for videos: Slowness, Shuffle&Learn, Order, Odd-One-Out, AoT, ST-Puzzle, DynamoNet, DPC, CBT, SpeedNet, MemDPC, CoCLR
- Datasets are balanced, so methods can take advantage of this

## Phase 2: the expansion phase

- **Applicable tasks**, beyond representation learning, including: standard computer visions tasks like tracking, localization, segmentation; few-shot learning
- **Multiple-modalities**: audio, video, text, ... more opportunities for supervision
- Training on **larger datasets**
- **More is better**: more data, longer training, more proxy tasks, more depth/width in the network
- **Datasets still tend to be curated**: AudioSet, IG65M, YouTube8M, HowTo100M, ...
- Good examples of exploring the benefits of more data:
  - Scaling and Benchmarking Self-Supervised Visual Representation Learning, Priya Goyal, Dhruv Mahajan, Abhinav Gupta, Ishan Misra, <https://arxiv.org/abs/1905.01235>
  - Evolving Losses for Unlabeled Video Representation Learning, AJ Piergiovanni, Anelia Angelova, Michael Ryoo, CVPR 2020

# Self-Supervised Learning



The Scientist in the Crib: What Early Learning Tells Us About the Mind  
by Alison Gopnik, Andrew N. Meltzoff and Patricia K. Kuhl

The Development of Embodied Cognition: Six Lessons from Babies  
by Linda Smith and Michael Gasser

# Phase 3: The Uncurated phase

- Self-supervision from uncurated data, i.e. no pre-defined datasets, instead:
  - Random YouTube videos, so not class balanced, long tailed
  - Daily life videos, e.g. Vlogs, babycams,
- New learning schedules:
  - Curriculum learning
  - How to obtain informative (hard) samples?
- More ambitious tasks ... discrete objects, memory
- Universal networks: able to ingest multiple-modalities and carry out multiple tasks
- Curated datasets still have their uses: become new evaluation benchmarks

# Summary

- Three phases of self-supervised learning
  - Classical
  - Expansion
  - Uncurated

Each stage has value for applications. Uncurated is less explored.

- Multiple-modality as free form of co-supervision in video
- Opportunity for learning more challenging applicable tasks