Vision and Sound

Computer Vision Fall 2018 Columbia University

Single-modality video representations





Slide credit: Andrew Owens



(McGurk 1976)

BBCTW

Same audio, different video!



(McGurk 1976)

Object Recognition







Natural Synchronization



Sound

Vision

Millions of Unlabeled Videos

rhythm

SoundNet





yard: 9.89%

lawn: 7.39%

English springer: 8.65% Welsh springer spaniel: 2.20% Border collie: 1.65%

restaurant: 11.79% dining room: 7.18% coffee shop: 6.70%

candle: 15.90% restaurant: 6.83% groom: 2.78%

Sound Recognition

Classifying sounds in ESC-50

Method	Accuracy	
Chance	2%	

Human Consistency

81%

Sound Recognition

Classifying sounds in ESC-50

Method	Accuracy
Chance	2%
SVM-MFCC	39%
Random Forest	44%
CNN, Piczak 2015	64%

Human Consistency

Sound Recognition

Classifying sounds in ESC-50

Method	Accuracy	
Chance	2%	
SVM-MFCC	39%	
Random Forest	44%	
CNN, Piczak 2015	64%	
SoundNet	74% 🗡 10% ga	IN
Human Consistency	81%	

Vision vs Sound

Low-dimensional embeddings via Maaten and Hinton, 2007



urban

nature

work-home

music–entertainment

sports

vehicles



Sound

Vision

Sensor Power Consumption







Layer 1











Man allen many with

~ ~ ~ ~ MAMMA MAMMA



Layer 5



Smacking-like

Layer 5



Chime-like







Scuba-like





Parents-like

Audiovisual Grounding



Which regions are making which sounds?



Audiovisual Grounding





Which objects make which sounds?



The sound of clicked object



The sound of clicked object



The sound of clicked object



Collect unlabeled videos









Mix Sound Tracks



How to recover originals?

Audio-only:

- Ill-posed
- permutation problem



Vision can help



Audiovisual Model


Audiovisual Model





Audiovisual Model



Original Audio



What does this sound like?



What does this sound like?



What does this sound like?



What regions are making sound?

Original Video













Estimated Volume

What sounds are they making?

Original Video



Embedding (projected and visualized as color)

Adjusting Volume







Volume 2	
•	×

Learning audio-visual correspondences



→ real or fake?

Learning audio-visual correspondences



Idea #1: random pairs



Arandjelovic, Zisserman. ICCV 2017

Audio-visual correspondence detector network



Vision hidden units



Sound hidden units



Sound Recognition

(a) ESC-50		(b) DCASE	
Method	Accuracy	Method	Accuracy
SVM-MFCC [26]	39.6%	RG [27]	69%
Autoencoder [2]	39.9%	LTT [19]	72%
Random Forest [26]	44.3%	RNH [28]	77%
Piczak ConvNet [25]	64.5%	Ensemble [32]	78%
SoundNet [2]	74.2%	SoundNet [2]	88%
Ours random	62.5%	Ours random	85%
Ours	79.3%	Ours	93%
Human perf. [26]	81.3%		

Visual Recognition

Method	Top 1 accuracy
Random	18.3%
Pathak <i>et al</i> . [24]	22.3%
Krähenbühl <i>et al</i> . [16]	24.5%
Donahue <i>et al</i> . [7]	31.0%
Doersch <i>et al</i> . [6]	31.7%
Zhang <i>et al</i> . [36] (init: [16])	32.6%
Noroozi and Favaro [21]	34.7%
Ours random	12.9%
Ours	32.3%

Linear classifier on top of features (ImageNet)

Idea #1: random pairs



Idea #2: time-shifted pairs





Idea #2: time-shifted pairs



Fused audio-visual representation



Fused audio-visual representation



What does the network learn?

Aligned vs. misaligned Aligned vs. misa

Class activation map (Zhou et al. 2016)

Top responses per category (speech examples omitted)

top baller

Dribbling basketball

ILLERBOOTCAMPLCOM

1.81

Dribbling basketball

Dribbling basketball





Playing organ



Chopping wood



Chopping wood

Application: on/off-screen source separation



Task: separate on-screen sounds from background noise

Creating training data



Skrthetensounofffesteren



On/off-screen source separation



On/off-screen source separation




Slide credit: Andrew Owens

Input video

.

OCBSN

On-screen prediction



Off-screen prediction



Input video

ONE-ON-ONE

TRUMP CALLS FOR DOJ INVESTIGATION OF NY TIMES OP-ED



Las Vegas

.....

6:04 PM PT

CUOMO PRIME TIME

On-screen prediction



TRUMP CALLS FOR DOJ INVESTIG

On-screen prediction

