# *Speech2Action:* Cross-modal Supervision for Action Recognition

Arsha Nagrani[1,2]   Chen Sun[2]   David Ross[2]   Rahul Sukthankar[2]   Cordelia Schmid[2]   Andrew Zisserman[1,3]

[1]VGG, Oxford   [2]Google Research   [3]DeepMind

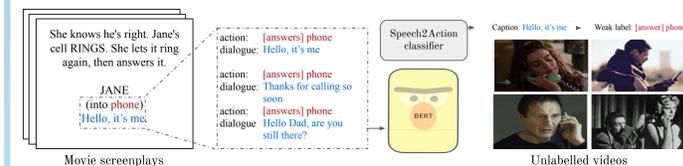CVPR SEATTLE
JUNE 14-19 2020 WASHINGTON

## Problem Definition and Contribution

**Goal:** Action recognition in video using only transcribed speech as supervision



### Motivation:
- Manual annotation of actions is expensive, not scalable
- Audiotrack is often freely available!

### Key Contributions:
- A Speech2Action model trained from screenplays that predicts actions from transcribed speech *alone*
- Applying Speech2Action to a unlabelled video, we obtain weak action labels for > 800K video clips
- An visual action classifier trained on these clips with *no* fine-tuning gets SOTA performance

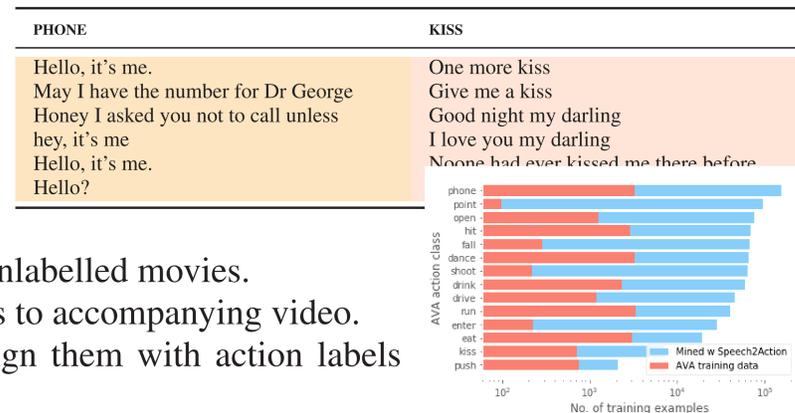## Mining with the Speech2Action Model

**Main idea:** Train a text-based model on movie scripts to predict actions from transcribed speech alone. Apply this to the transcribed speech from unlabelled videos to automatically get weak labels for video.

### Speech2Action Model
- Finetune a pretrained BERT model on speech-action pairs in moviescripts

### Mining Clips Automatically:
- Apply Speech2Action to subtitles of unlabelled movies.
- Assign label for highly confident preds to accompanying video.
- Mine over 800K video clips and assign them with action labels based on the speech alone.

## IMSDb Dataset

- 1,080 movie scripts from www.IMSDb.com, 22 genres
- Create a text dataset of **speech** paired with **action** labels from the scene directions, based on proximity in the movie script

### Examples of Movie Scripts



## Results on Visual Action Recognition

### Examples of clips mined using Speech2Action:



### Examples of abstract actions mined using Speech2Action:



### Results on 14 AVA mid and tail classes

| Data | Per-Class AP | | | | | | | | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | drive | phone | kiss | dance | eat | drink | run | point | open | hit | shoot | push | hug | enter |
| AVA (fully supervised) | 0.63 | 0.54 | 0.22 | 0.46 | 0.67 | 0.27 | 0.66 | 0.02 | 0.49 | 0.62 | 0.08 | 0.09 | 0.29 | 0.14 |
| S2A-mined (zero-shot) | 0.83 | 0.79 | 0.13 | 0.55 | 0.68 | 0.30 | 0.63 | 0.04 | 0.52 | 0.54 | 0.18 | 0.04 | 0.07 | 0.04 |
| S2A-mined + AVA | **0.86** | **0.89** | **0.34** | **0.58** | **0.78** | **0.42** | **0.75** | 0.03 | **0.65** | **0.72** | **0.26** | **0.13** | **0.36** | **0.16** |

### Visual Action Recognition
- Train S3D-G with an 18-way softmax
- Evaluate on AVA with NO finetuning, on mid and tail classes. These actions occur *rarely*. For 8 classes, exceed fully supervised performance without any supervision.
- On HMDB51, obtain 17% improvement over training from scratch and outperform SOTA.
- Even able to label abstract actions like count, follow etc!

### Results on HMDB51

| Method | Architecture | Pre-training | Acc. |
| --- | --- | --- | --- |
| Shuffle&Learn | S3D-G (RGB) | UCF101 | 35.8 |
| OPN | VGG-M-2048 | UCF101 | 23.8 |
| ClipOrder | R(2+1)D | UCF101 | 30.9 |
| 3DRotNet | S3D-G (RGB) | Kinetics | 40.0 |
| DPC | 3DResNet18 | Kinetics | 35.7 |
| CBT | S3D-G (RGB) | Kinetics | 44.6 |
| DisInit (RGB) 2019 | R(2+1)D-18 | Kinetics** | 54.8 |
| Korbar et al. 2018 | I3D (RGB) | Kinetics | 53.0 |
| - | S3D-G (RGB) | Scratch | 41.2 |
| Ours | S3D-G (RGB) | S2A-mined | **58.1** |

## Acknowledgments