# Audio-Visual Learning

Chuang Gan
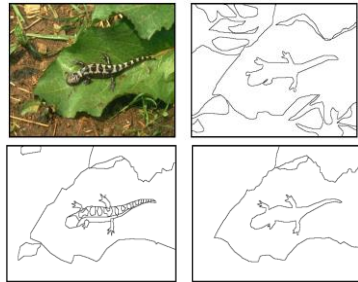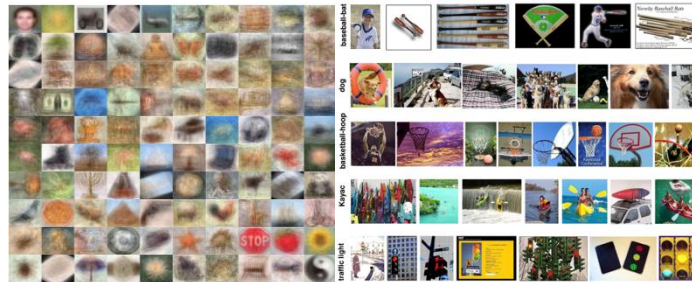
# Welcome Visitors!

# The Mcgurk Effect

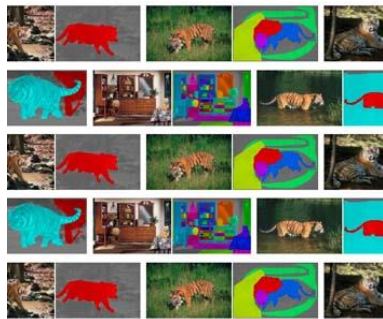# Learning from images and video frames


BSD (2001)


Caltech 101 (2004), Caltech 256 (2006)


PASCAL (2007-12)


LabelMe (2007)


ImageNet (2009)


SUN (2010)


UCF-101 (2012)


Youtube-8M (2017)


Kinetics (2017)

4

# What can sound give us?

physical interactions

speech

sound of distant object

# Can machines connect sight with sound for rich perception?

# Task: Visual Sound Separation

Given a music performance video…

**Mixed sound**

# Task: Visual Sound Separation

…we aim to separate two sounds played by different instruments.
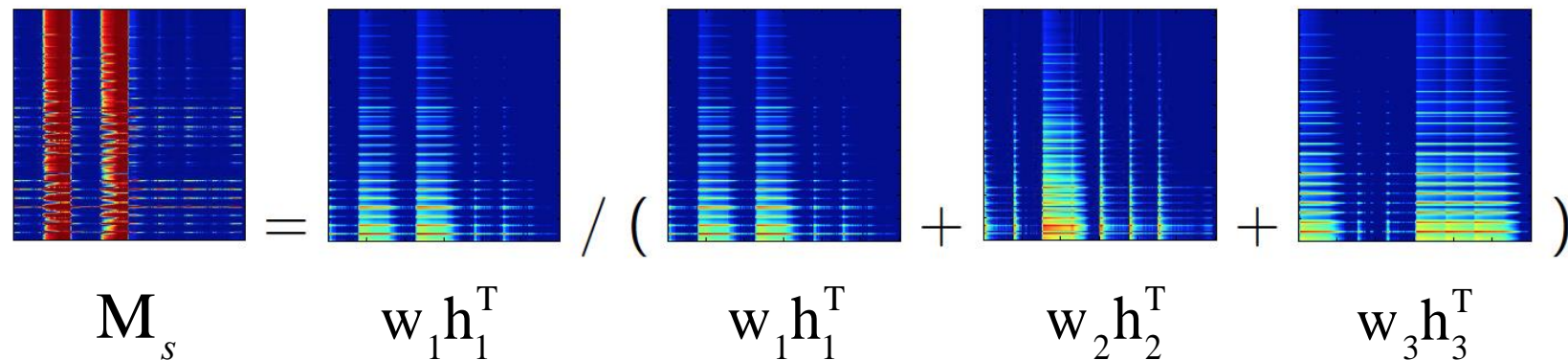
**Separated sound**

**Mixed sound**



**Network**

# Source Separation: Traditional Approach

❑To separate one component out of K:
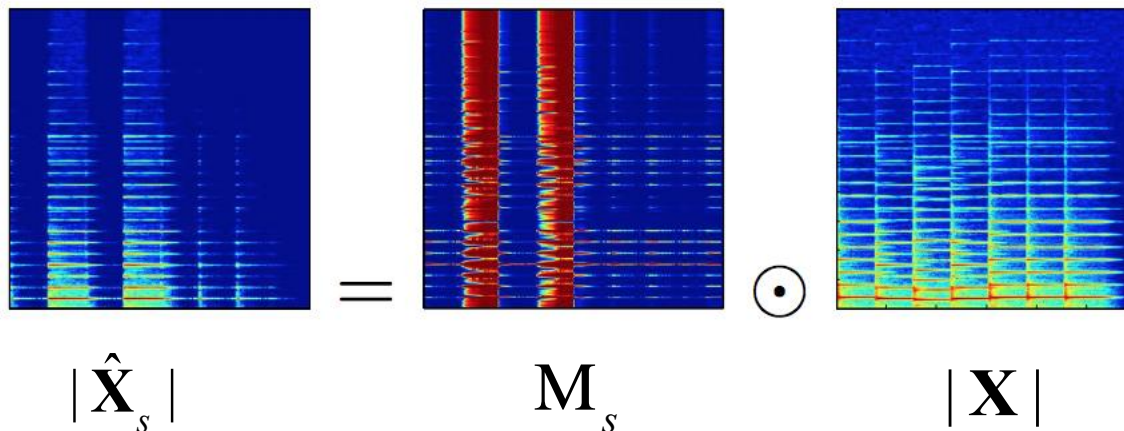- 1. Estimate the mask for the target component

$$\mathbf{M}_s = \frac{\mathbf{w}_1\,\mathbf{h}_1^{\mathrm{T}}}{\sum\limits_{i=1}^{K} \mathbf{w}_i\,\mathbf{h}_i^{\mathrm{T}}}$$



$$\mathbf{M}_s = \mathbf{w}_1\mathbf{h}_1^{\mathrm{T}} / ( \mathbf{w}_1\mathbf{h}_1^{\mathrm{T}} + \mathbf{w}_2\mathbf{h}_2^{\mathrm{T}} + \mathbf{w}_3\mathbf{h}_3^{\mathrm{T}} )$$

# Source Separation: Traditional Approach

❑To separate one component out of K:
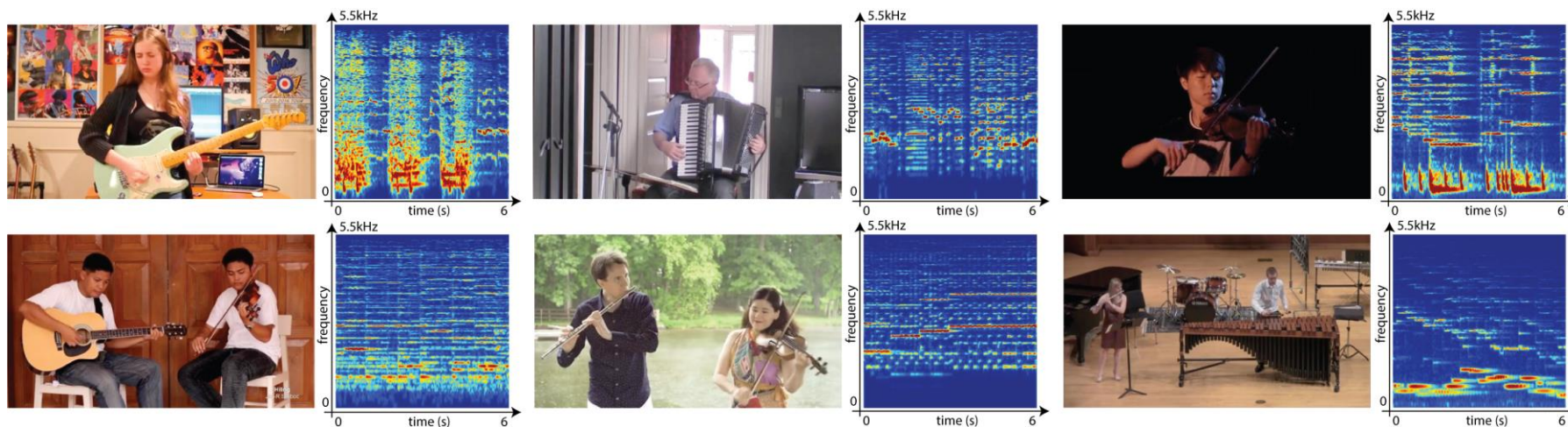- 2. Masking on the input spectrogram to separate the component

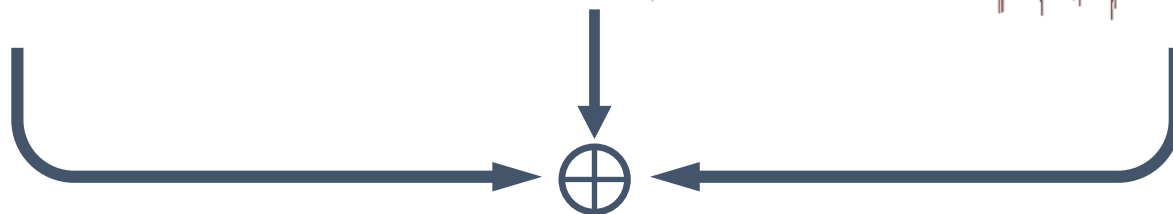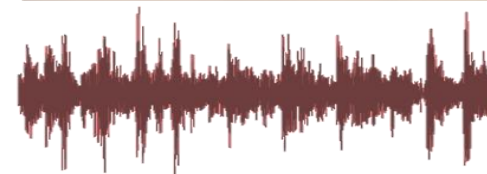$$|\hat{\mathbf{X}}_s| = \mathbf{M}_s \odot |\mathbf{X}|$$



$$|\hat{\mathbf{X}}_s| \quad = \quad \mathbf{M}_s \quad \odot \quad |\mathbf{X}|$$

# Our Ideas: Learning from Music Videos

❑ Internet music videos
- Keyword search without labeling
- >20 kinds of commonly seen musical instruments
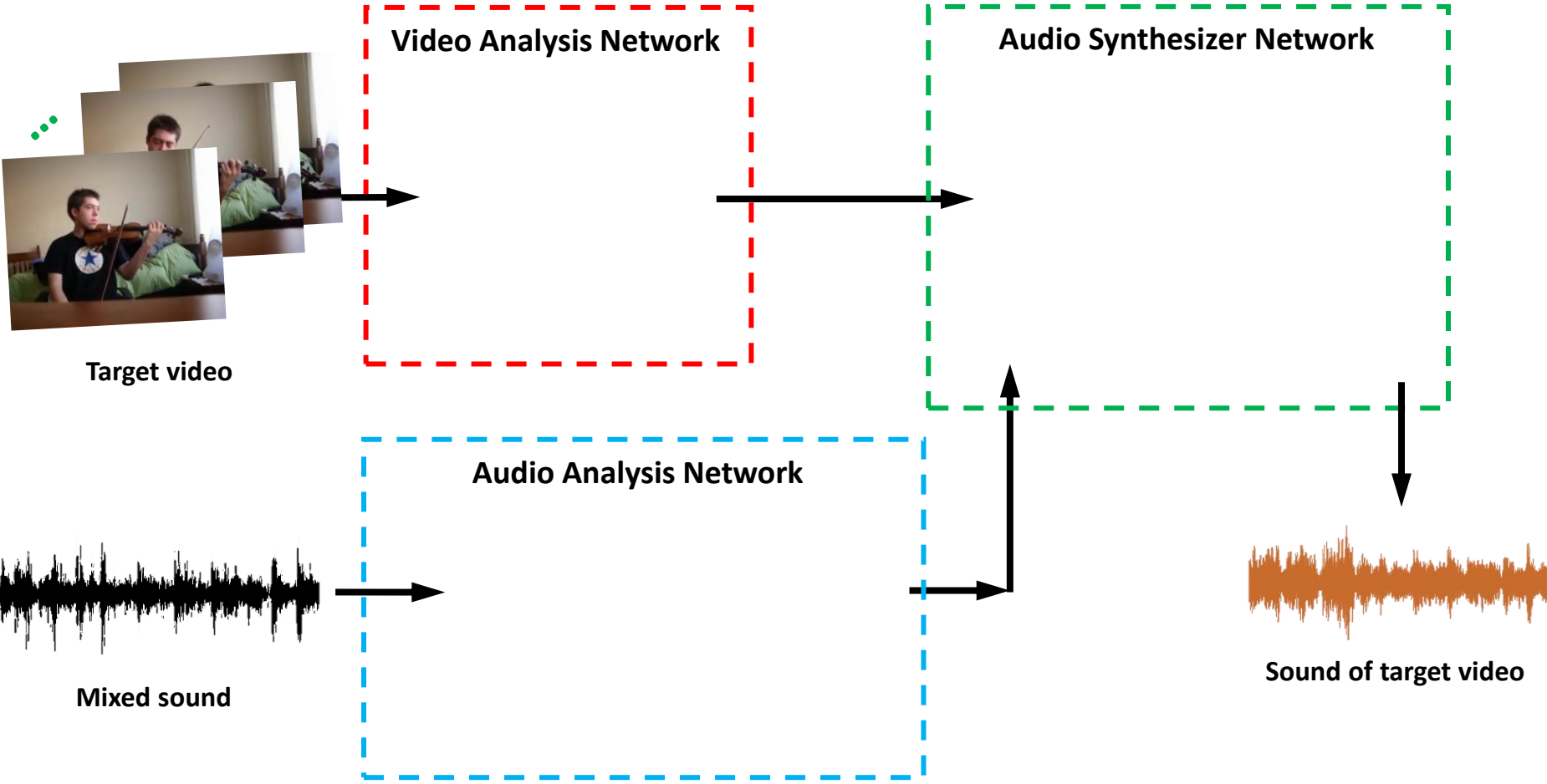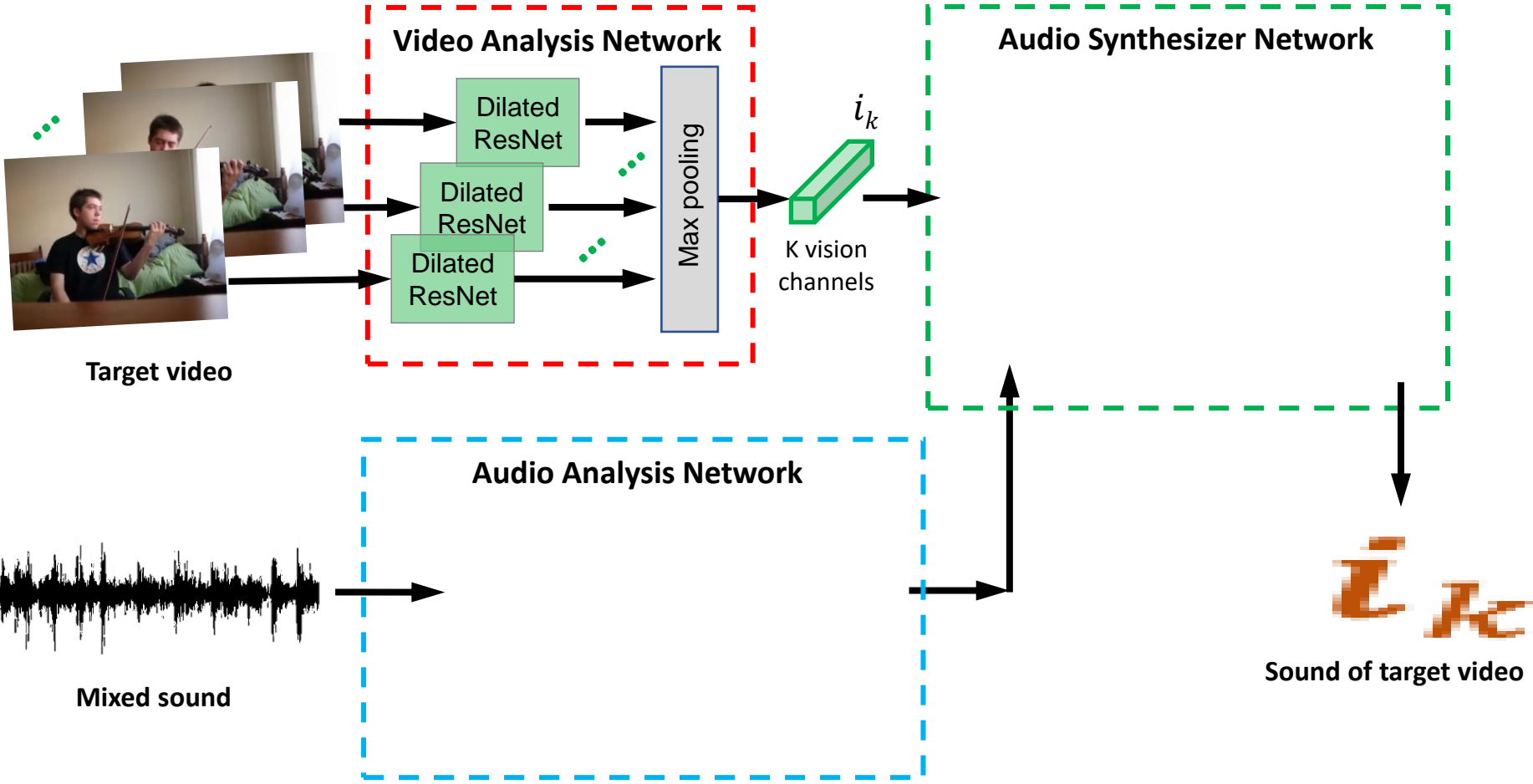- >1000 solos and duets



*Zhao, Gan et al. "The Sound of Pixels." ECCV 2020.*

# Mixing the Sound



*Zhao, Gan et al. "The Sound of Pixels." ECCV 2020.*

# Vision to Rescue for Self-supervised Prediction



**Visual-audio source separation model**

Target video

Video Analysis Network

Audio Synthesizer Network

Mixed sound

Audio Analysis Network

Sound of target video

# Mix-and-Separate Framework



Video Analysis Network

Audio Synthesizer Network

Target video

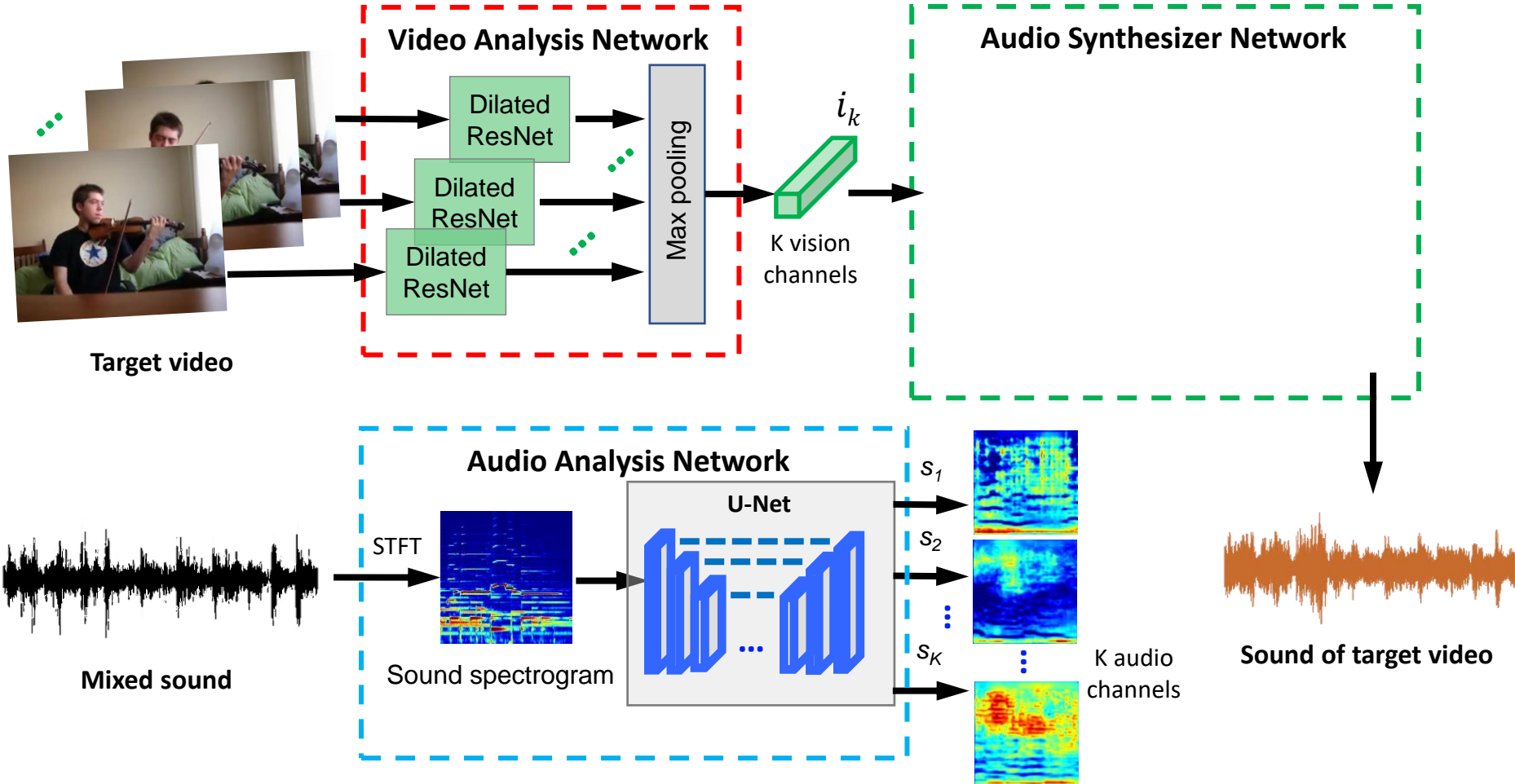Audio Analysis Network

Mixed sound
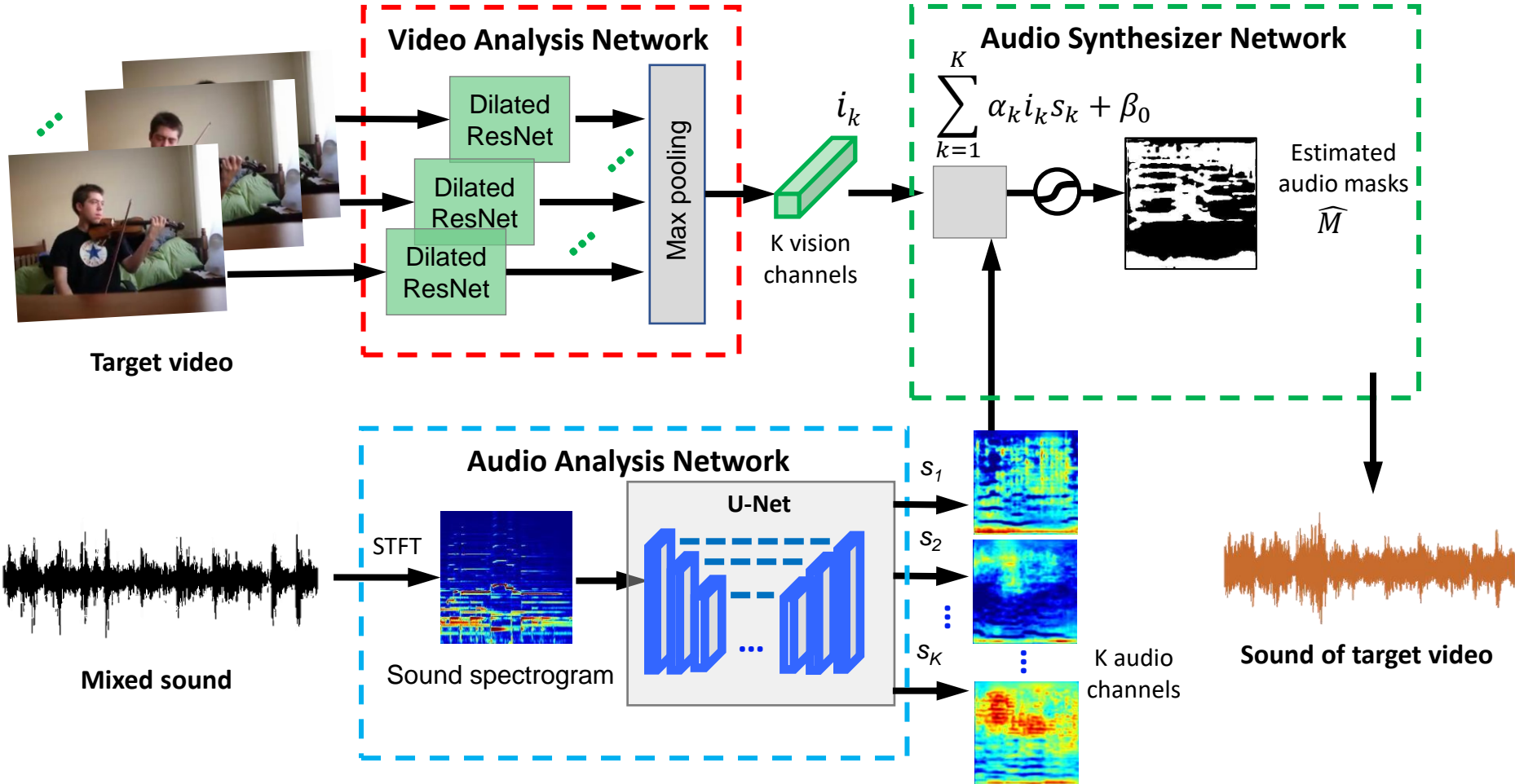
Sound of target video

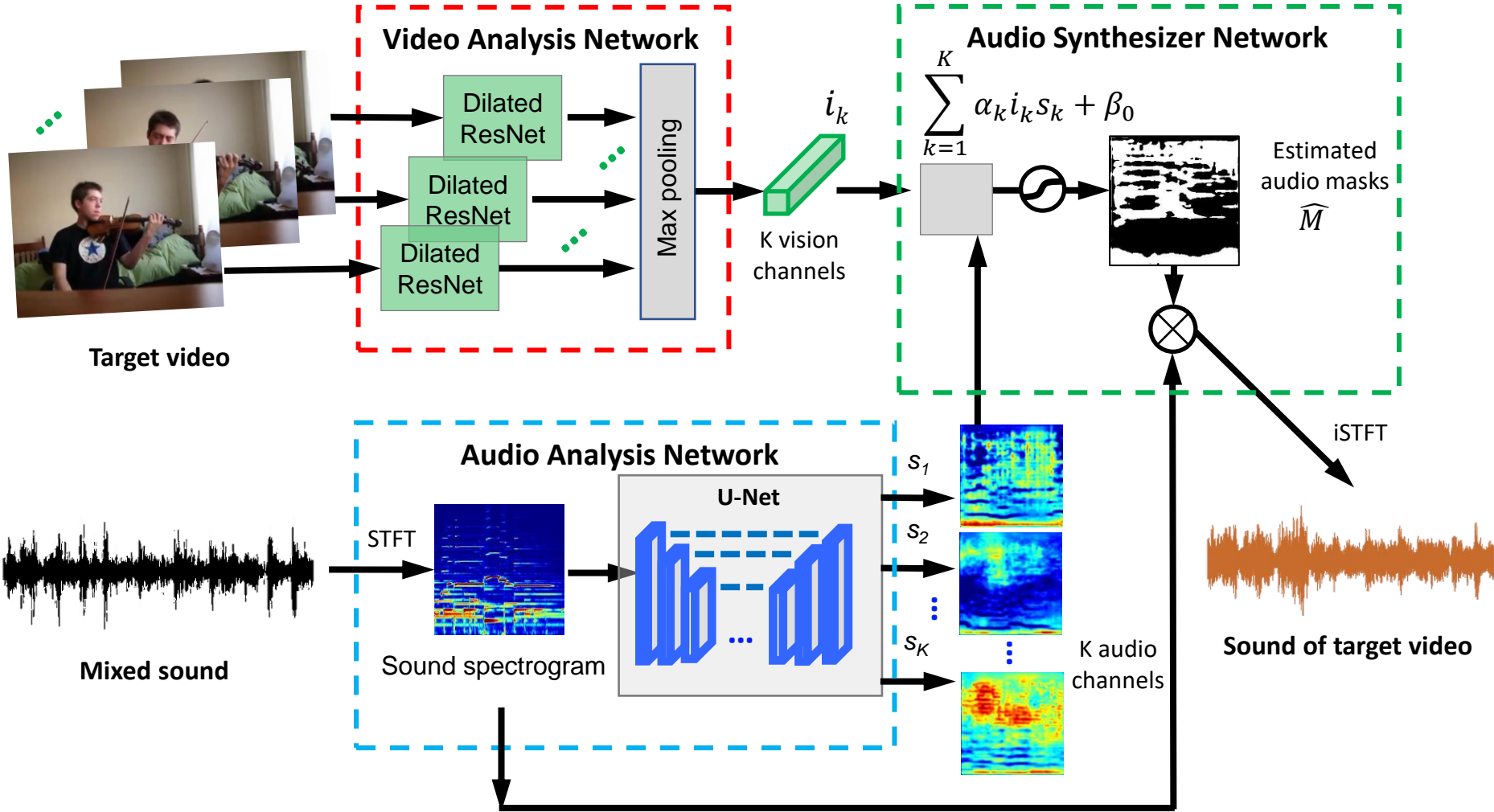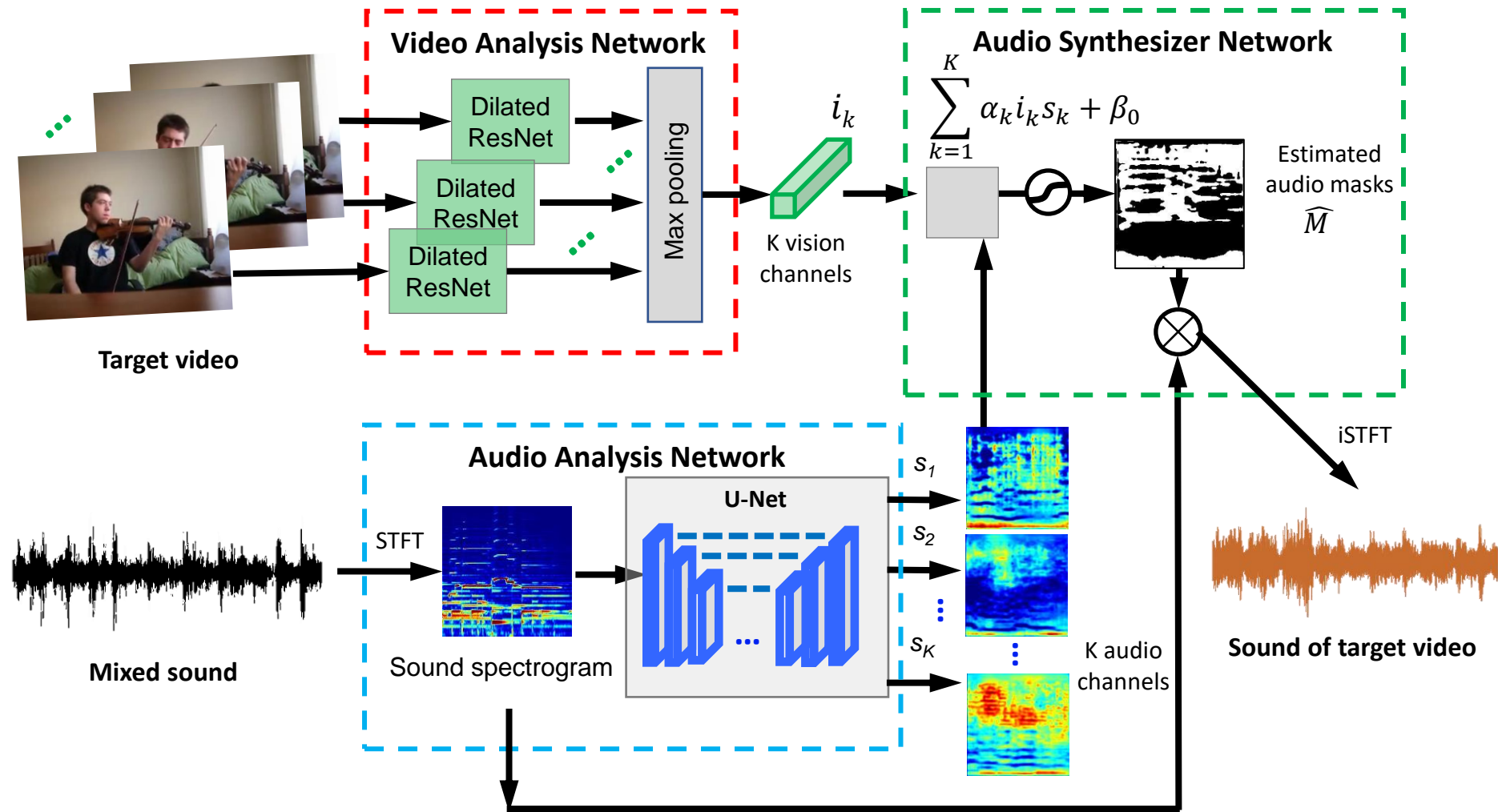# Mix-and-Separate Framework

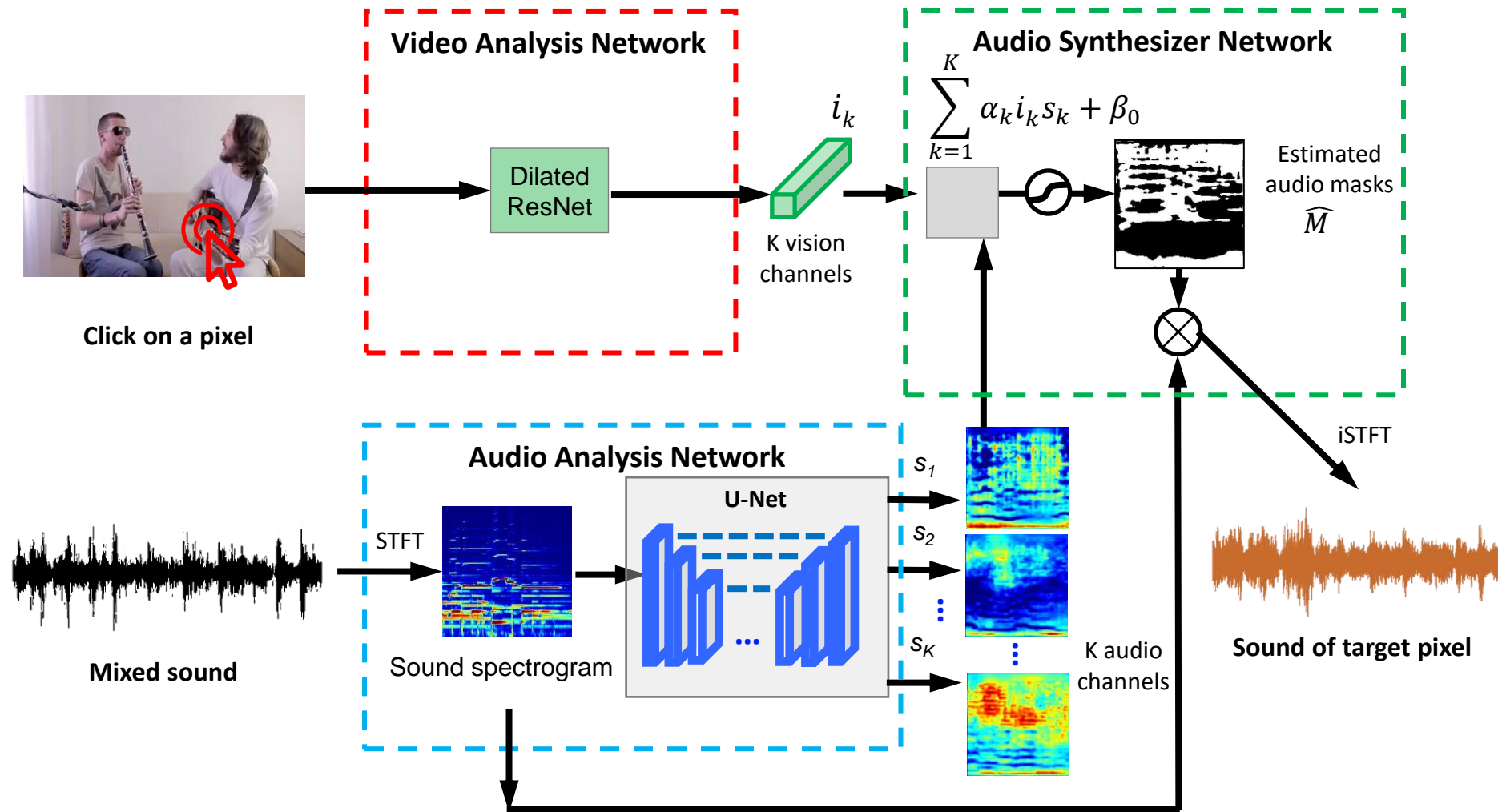# Mix-and-Separate Framework

# Mix-and-Separate Framework

# Mix-and-Separate Framework

# Test Time

# Test Time: using Pixel Feature instead

# Original Video

# The sound of clicked object…

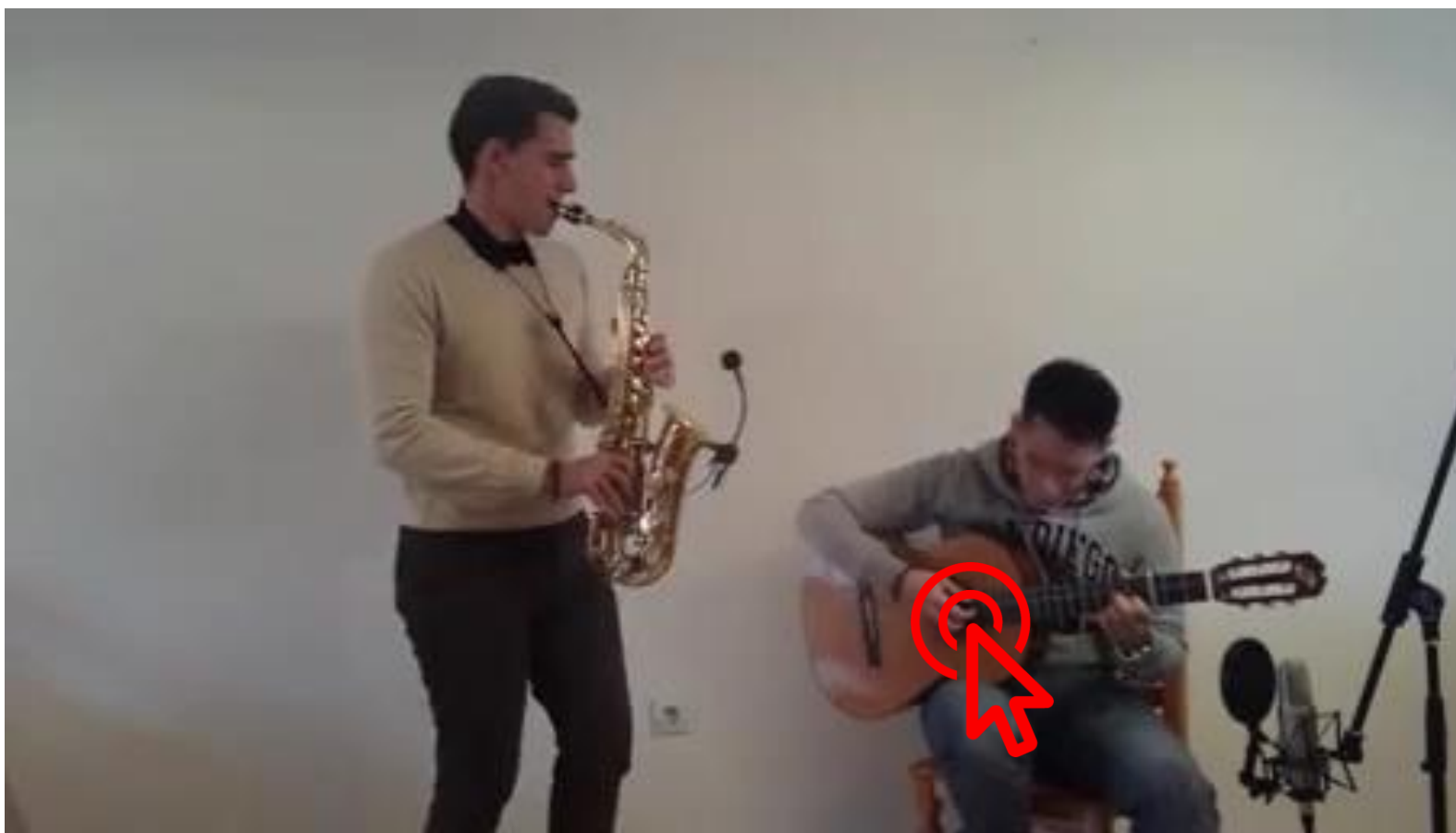# The sound of clicked object…
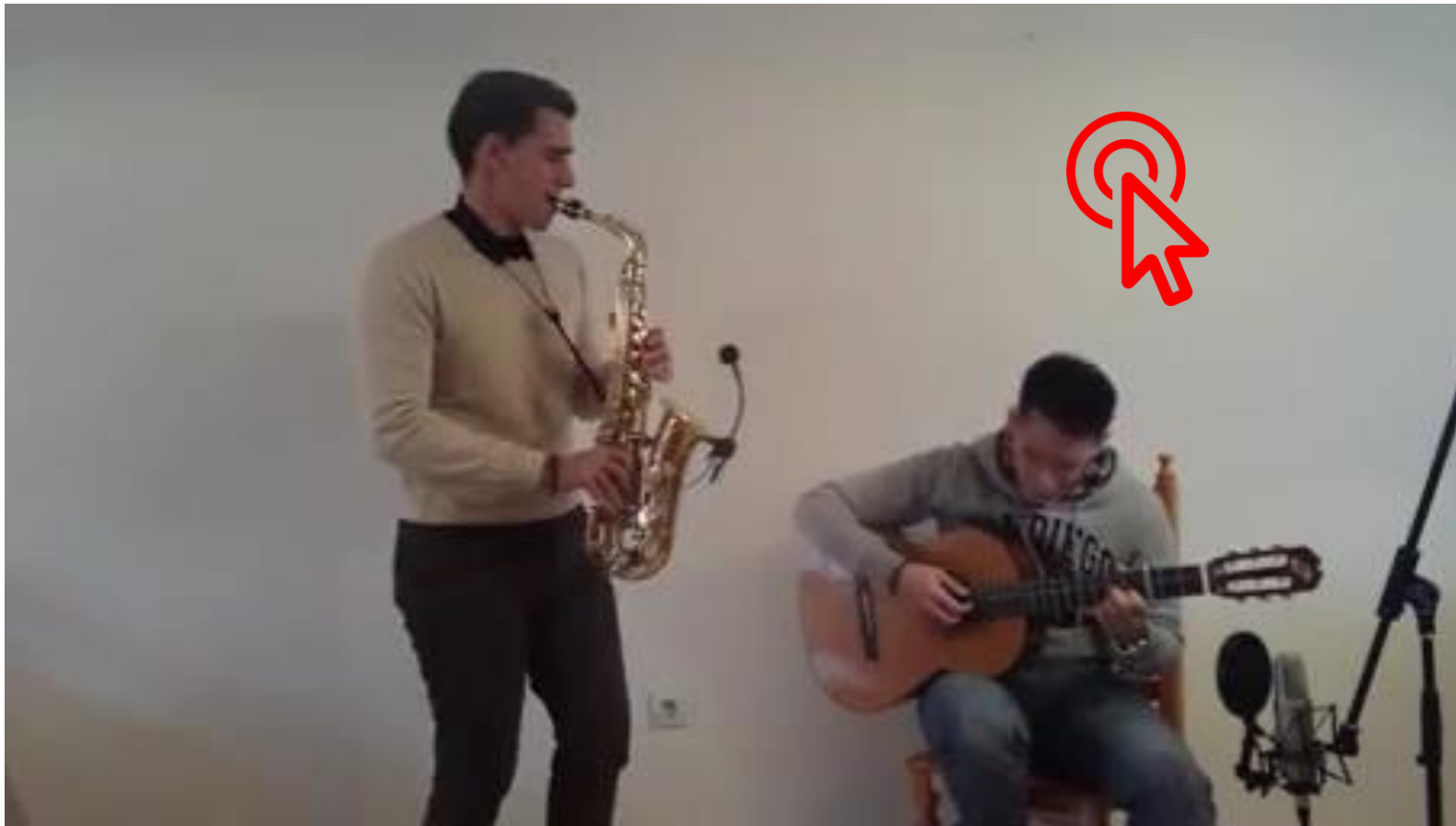
# The sound of clicked object…

# Original Video

# The sound of clicked object…

# The sound of clicked object…

# The sound of clicked object…

# Application: Music Remix

# Application: Music Remix

# Limitation

❑ Most existing methods use **raw pixel** or **optical flow** as input.

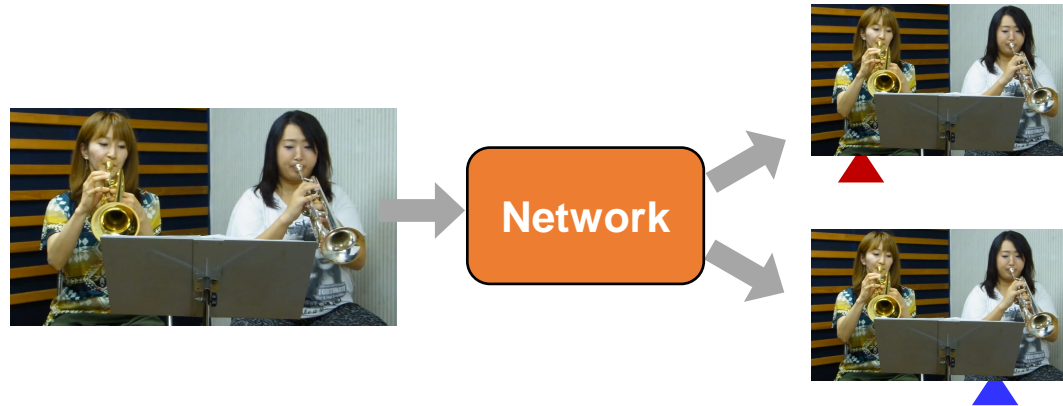❑ Problem: limited to <span style="color:red">**separate multiple instruments of the same types**</span>.
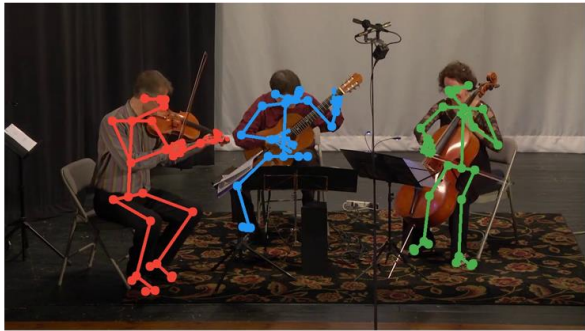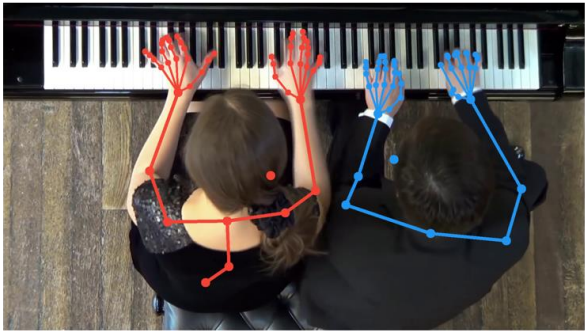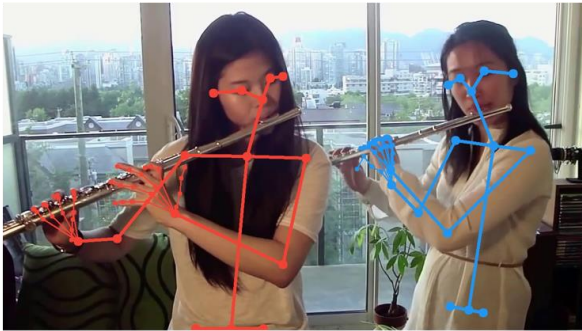
# Our Ideas

□ Problem: limited to **separate multiple instruments of the same types**.



□ We propose to Identifying a melody by studying a musician's body language using ``**Music Gesture**''.

*Gan et al. "Music Gesture for Visual Sound Separation." CVPR 2020.*

# Music Gesture

❑ Keypoint-based structured representations

# Visual sound separation results

# Sound of Motion



**Mixed sound**



**Separated sound1**



**Separated sound2**

# Music Gesture



**Mixed sound**

**Separated sound1**

**Separated sound2**

# Sound of Motion



**Mixed sound**

**Separated sound1**
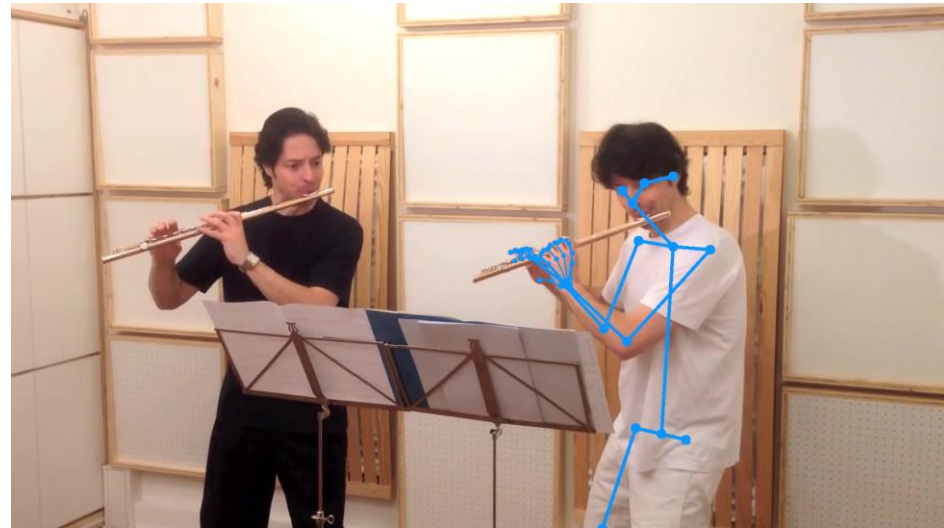
**Separated sound2**

# Music Gesture



Mixed sound

Separated sound1

Separated sound2

# Sound of Motion



**Mixed sound**
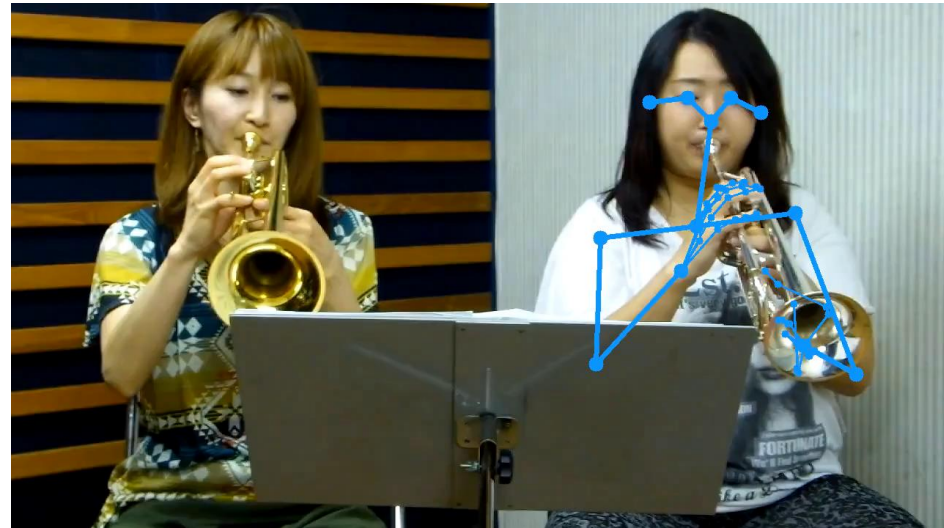


**Separated sound1**



**Separated sound2**

# Music Gesture



**Mixed sound**

**Separated sound1**

**Separated sound2**

# Multiple instruments

# Music Gesture



**Mixed sound**

**Separated sound1**

**Separated sound2**

# Music Gesture



Mixed sound



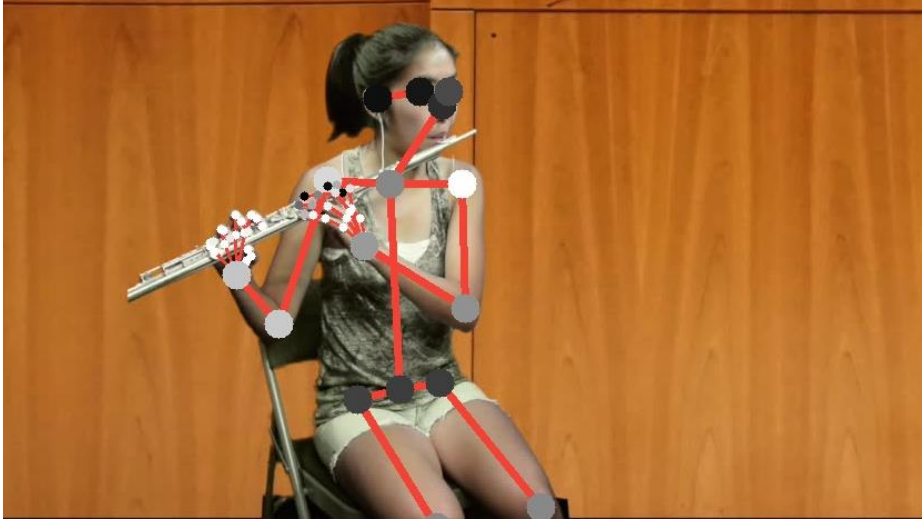Separated sound3



Separated sound4

# Attention Map of Key points

# The sound of body parts



Mixed sound

Separated sound1

Separated sound2

# Can we generate music from videos?

Given a silent music performance video…



**Silent music performance video**

*Gan et al. "Foley Music: Learning to Generate Music from Videos." ECCV 2020.*

# Can we generate music from videos?
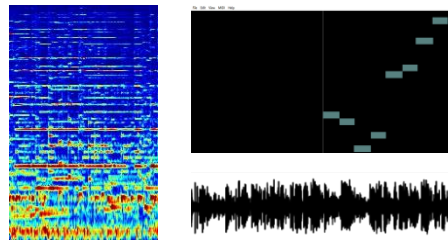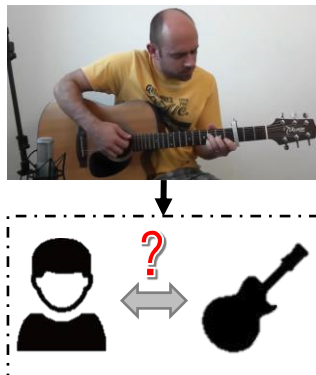
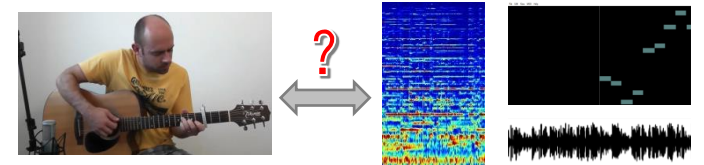…we aim to generate plausible music.



**Deep Neural Network**

**Silent music performance video**          **Performance with generated sound**

*Gan et al. "Foley Music: Learning to Generate Music from Videos." ECCV 2020.*

# Challenges

☐ Hard to learn visual-audio mappings from unlabeled video

☐ Three things matter：

◆ Visual perception module → interactions between instrument and player

◆ Audio representation → musical rules, easy to predict from visual signals

◆ Visual-audio model → association between two modalities

Choose ?

# Challenges

☐ Hard to learn visual-audio mappings from unlabeled video

☐ Three things matter：

- ◆ Visual perception module → interactions between instrument and player

- ◆ Audio representation → musical rules, easy to predict from visual signals

- ◆ Visual-audio model → association between two modalities

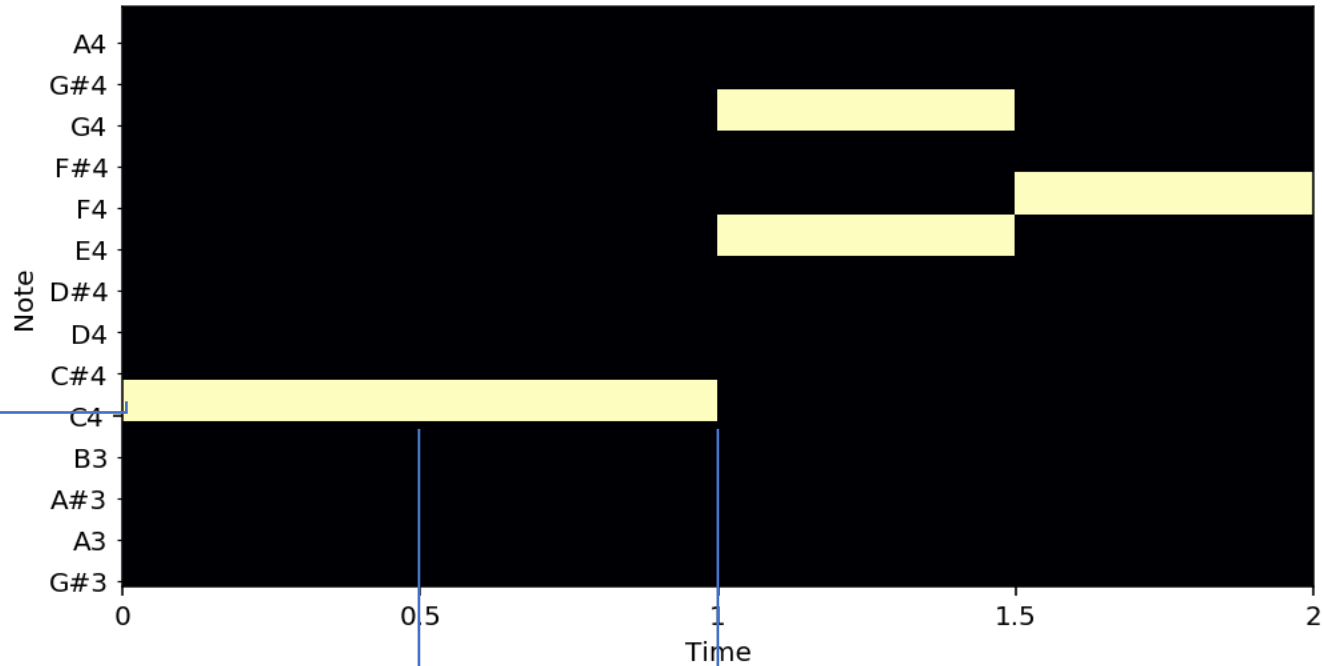We use **body keypoints** to explicitly model the body and finger.

We use **Musical Instrument Digital Interface (MIDI)** to represent music.

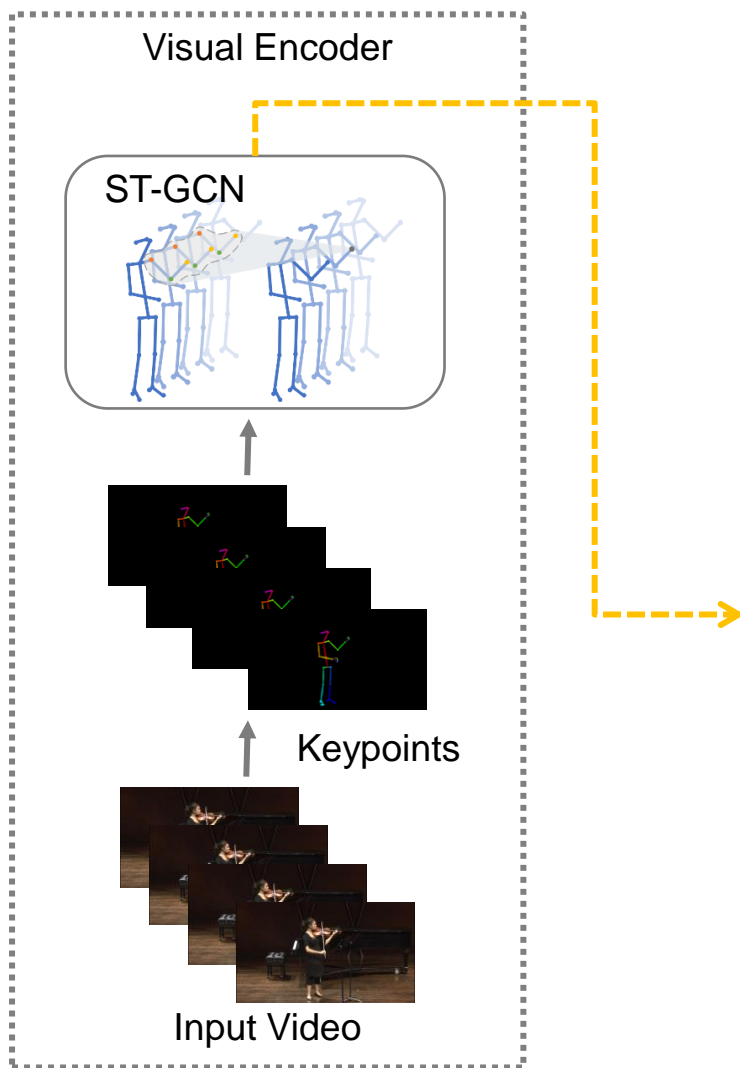# MIDI Event Representations



Velocity Event, Note On Event    Time Shift Event    Note Off Event

201          39          237          127

Each note represented as a sequence of MIDI Events

- ❏ Note On Event $\in [0, 88]$, based on Pitch
- ❏ Note Off Event $\in [88, 176]$, based on Pitch
- ❏ Velocity Event $\in [176, 208]$, based on Velocity
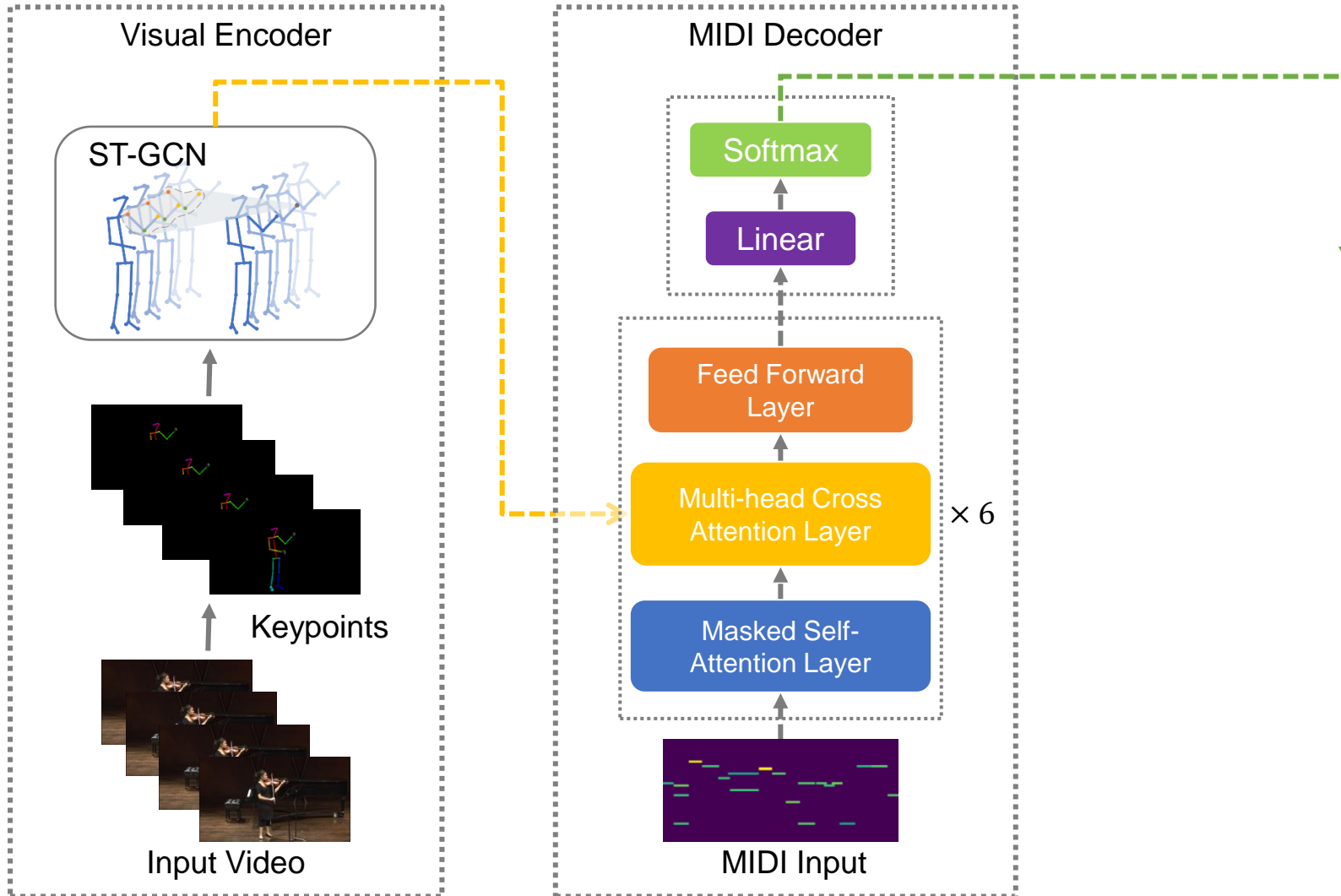- ❏ Time Shift Event $\in [208, 240]$, based on Duration

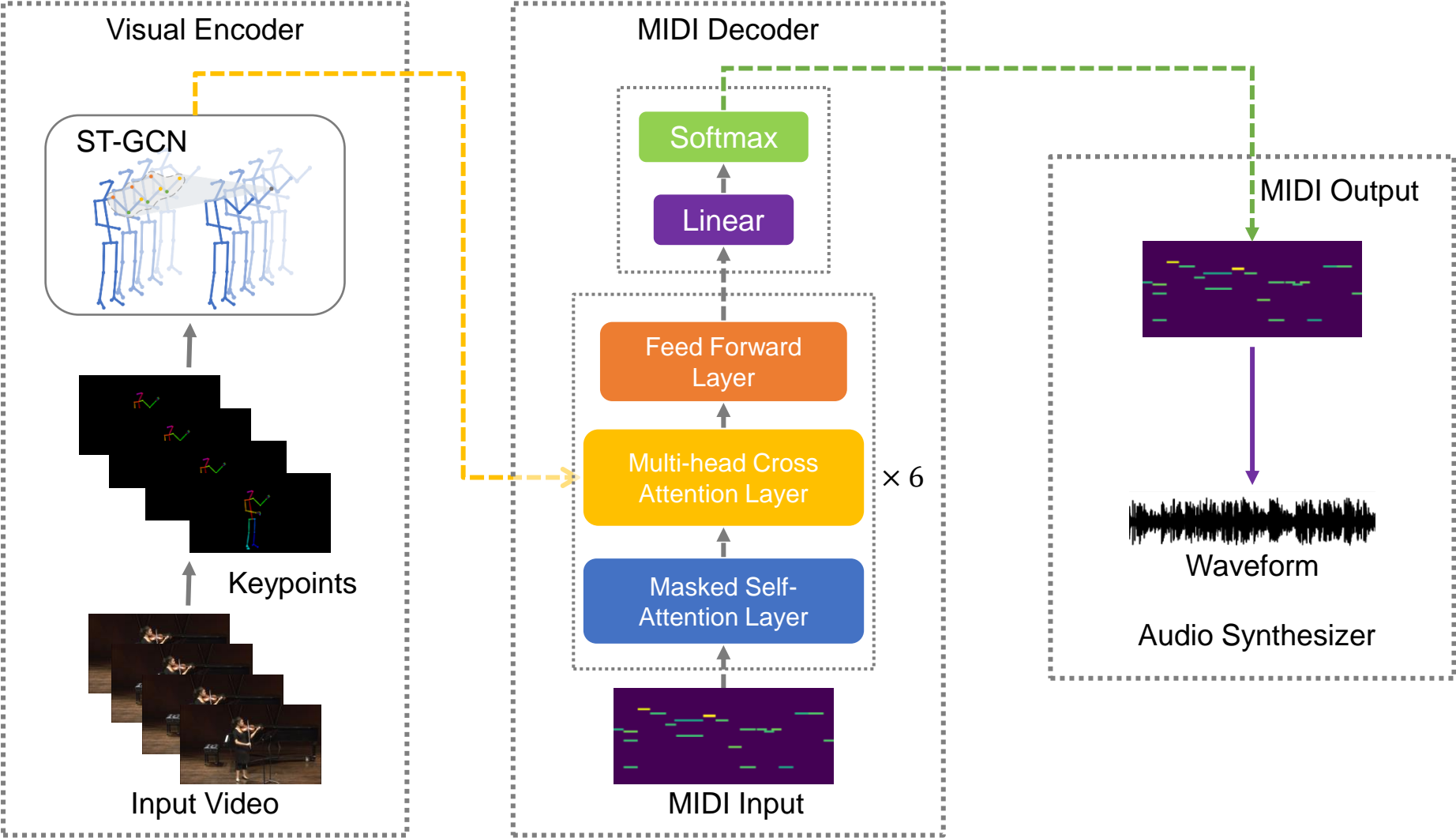# Our method

We first encode human keypoints by using graph CNN.



Visual Encoder

ST-GCN

Keypoints

Input Video

# Our method

We then translate the keypoints features into MIDI using Transformers.

# Our method

Finally, we synthesize audio from MIDI.

# Music generation results

# Ukulele

Piano

# Guitar

# Style editing results

# Bass



Original prediction

Style editing

A major

F major

G major

# Tuba



**Style editing**
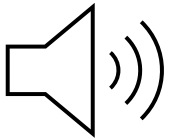
**A major**

**F major**
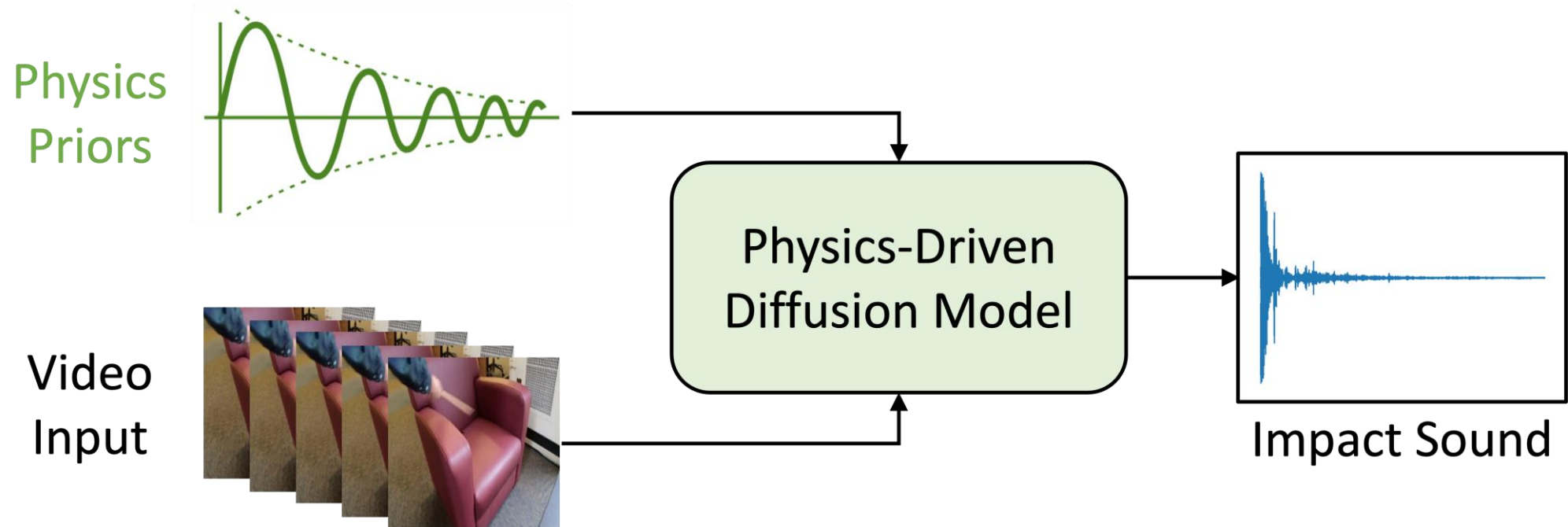
**G major**

**Original prediction**

An impact sound of physical object interactions is critical to perception
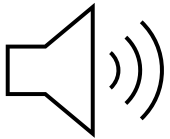
Could we generate the impact sound from vision?

# Our Framework



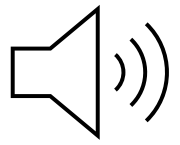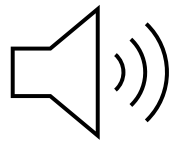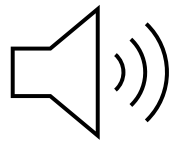Physics Priors
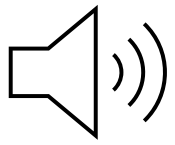
Video Input

Physics-Driven Diffusion Model

Impact Sound

Our model is applicable to
a variety of videos and materials

The transparency of **physics priors** allows us to perform interesting **sound editing**

# Physics Priors of **Glass** + Video Input



Original Result

# Physics Priors of **Glass** + Video Input



Original Result

Transformed Result

# Physics Priors of **Cloth** + Video Input



Original Result
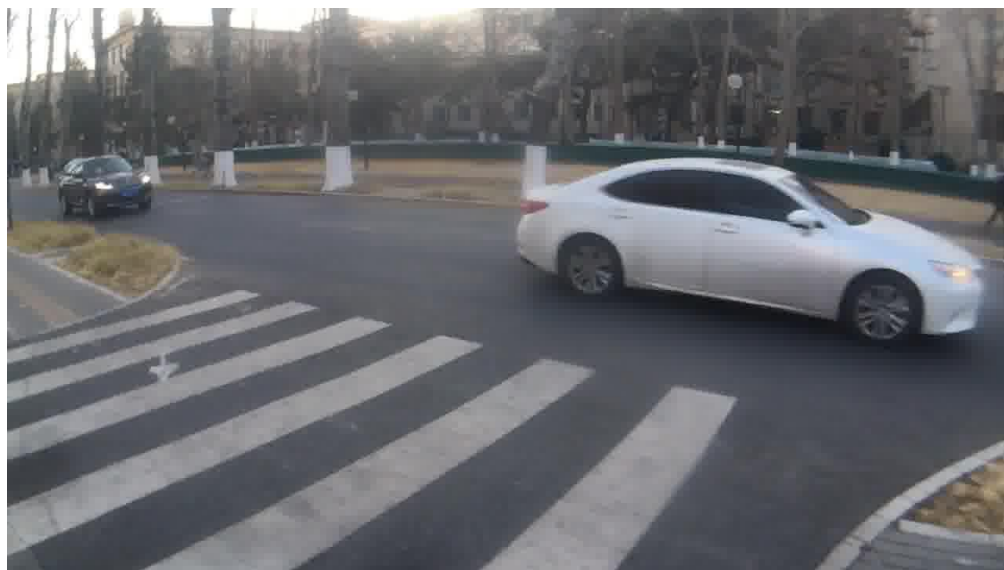
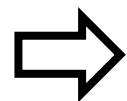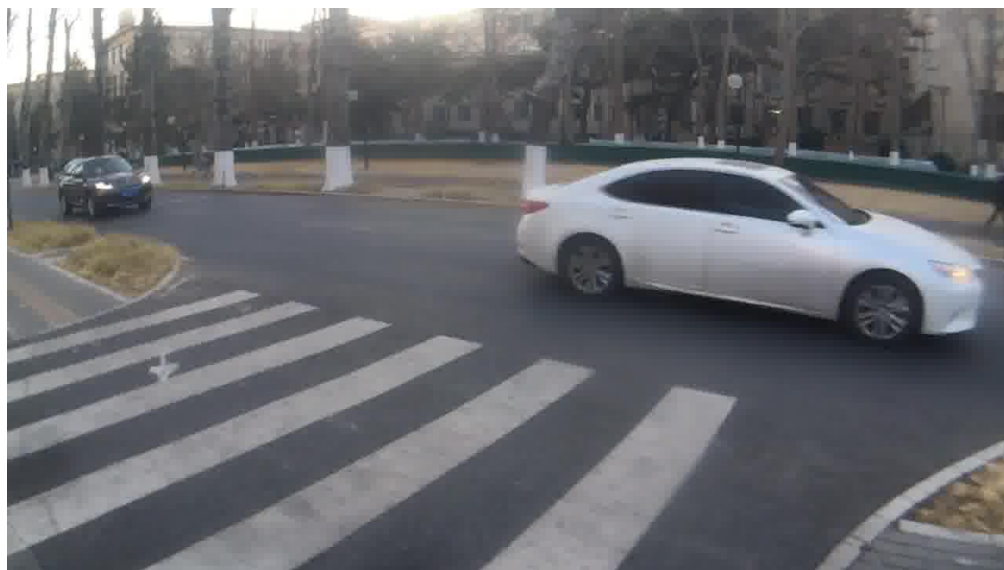# Physics Priors of **Cloth** + Video Input
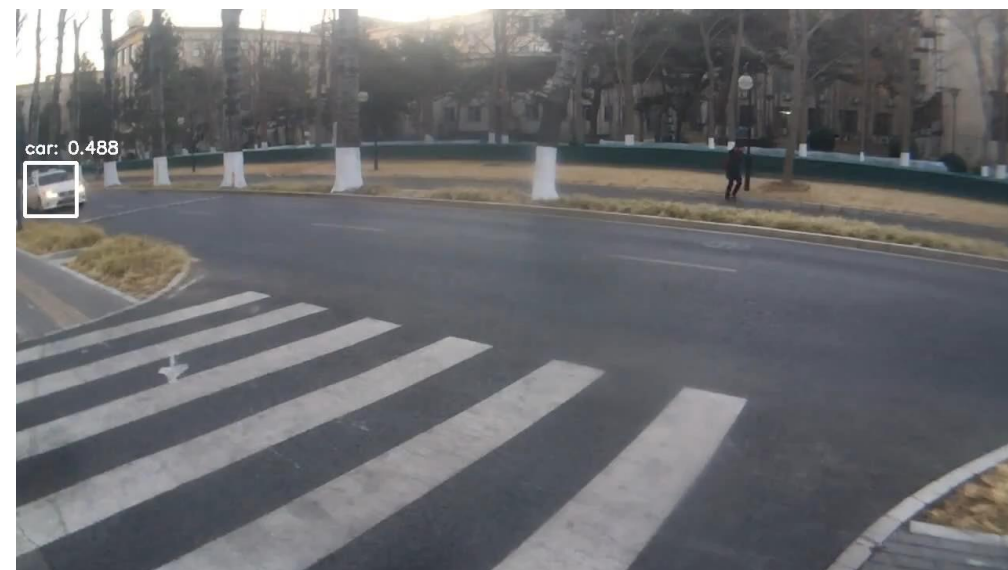


Original Result

Transformed Result

**Given an input video…**



Self-supervised moving vehicle tracking with stereo sound. Gan et.at. ICCV 19

**Given an input video…**
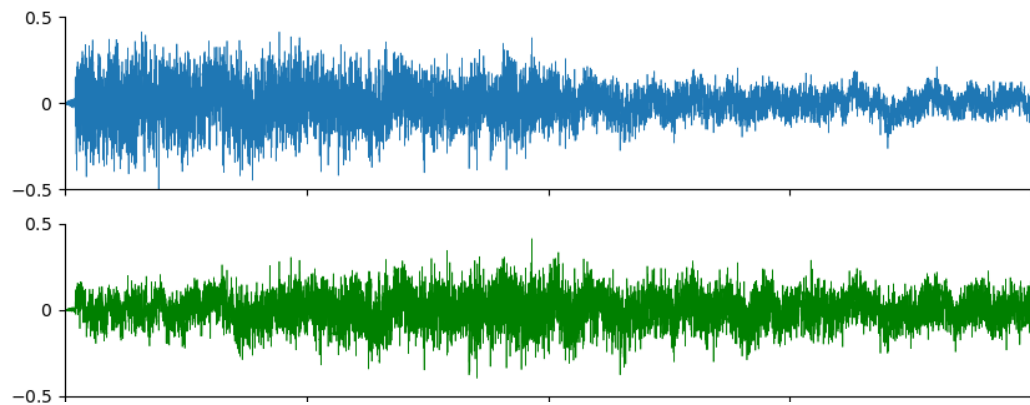
**Visual detection network can track the vehicles using visual input.**
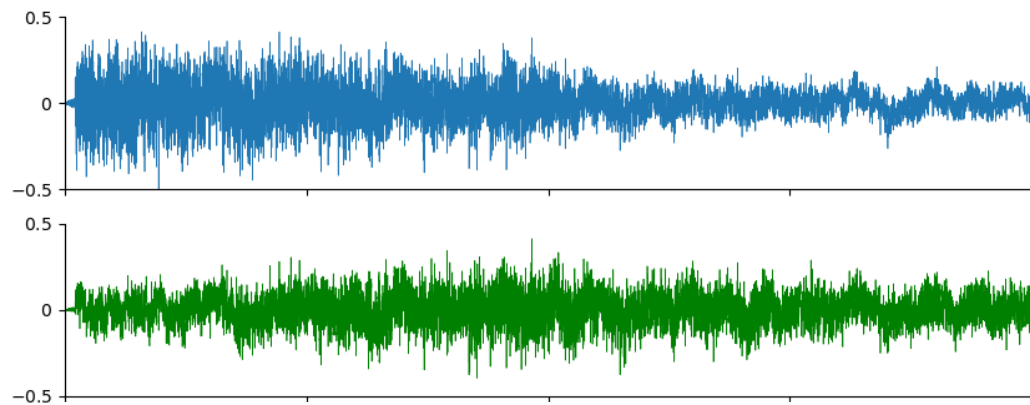


car: 0.488

**Visual tracking**

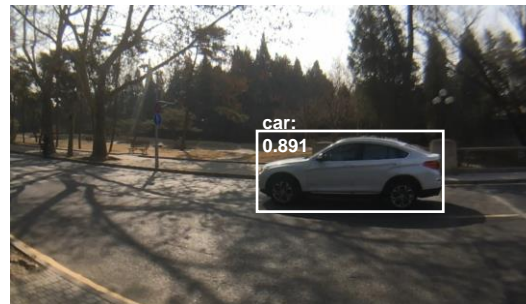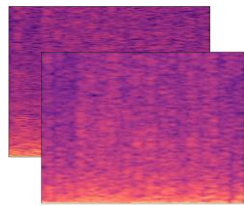# What if given a piece of input stereo sound only?

**What if given a piece of input stereo sound only?**



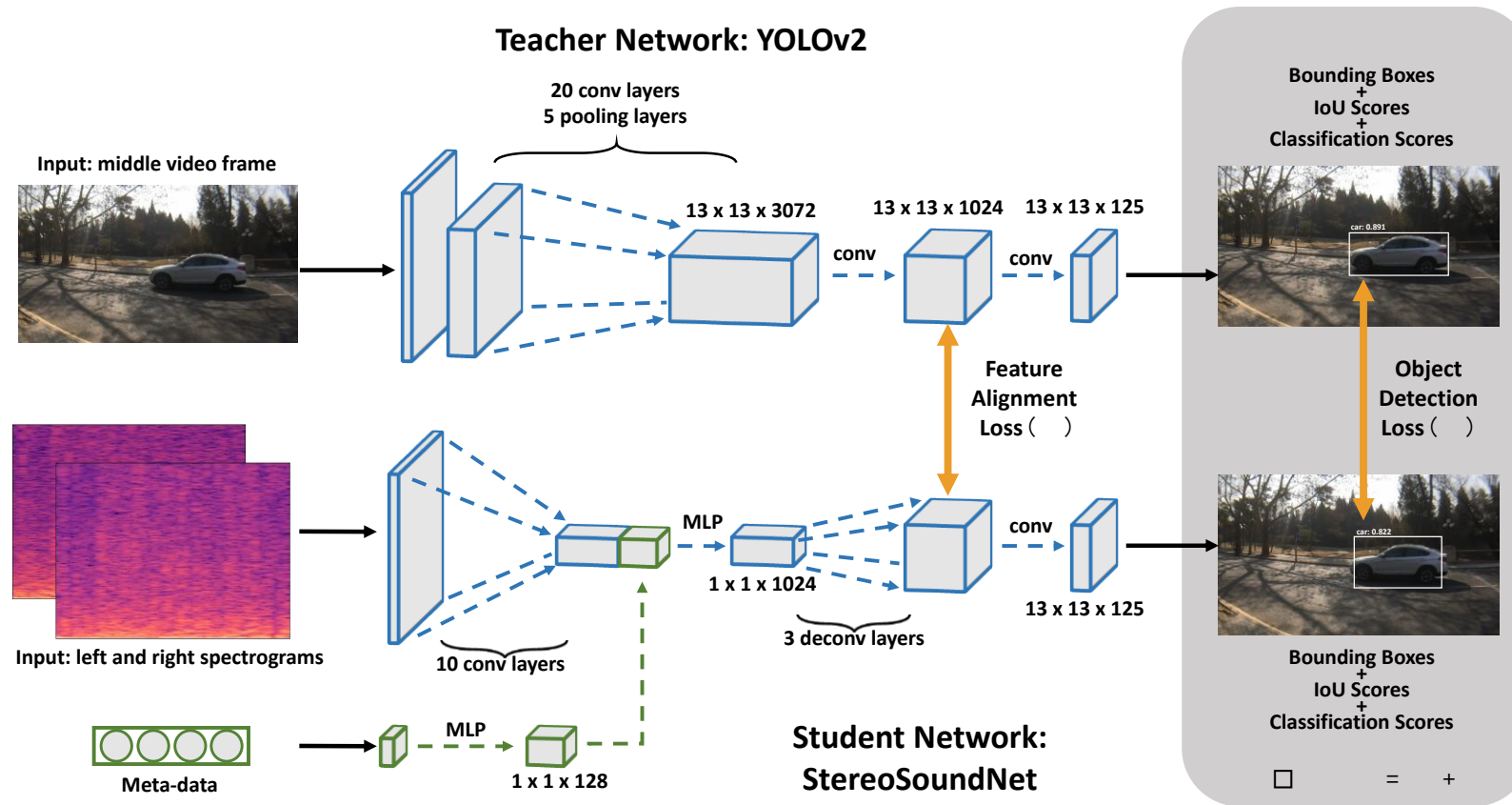**Where are the vehicles?**

# Applications

- Tracking under poor lighting scene
- Tracking under visual occlusion scene
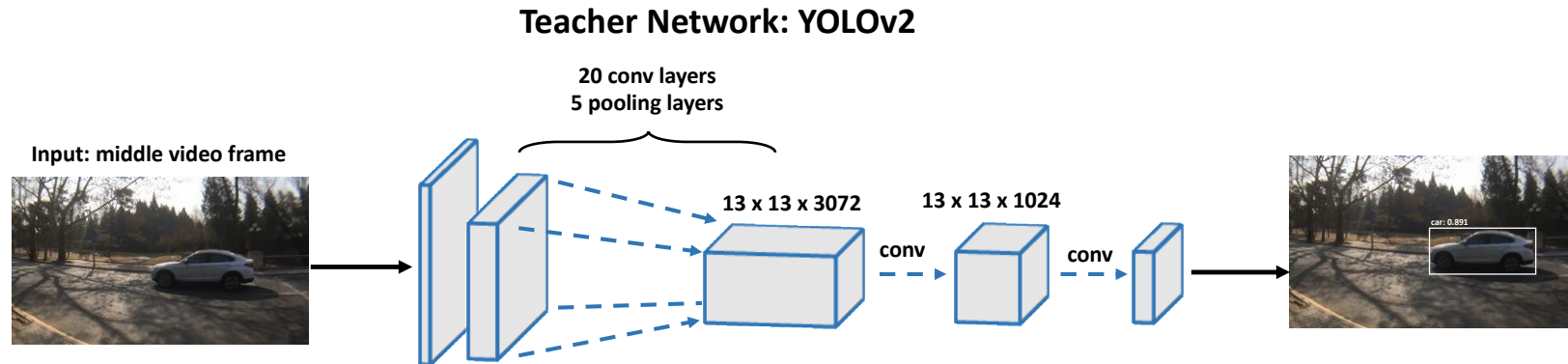- Energy-efficiency, privacy-preserving

# Our Methods

- Teacher-student alignment



**Teacher Network: YOLOv2**

Input: middle video frame

20 conv layers
5 pooling layers

13 x 13 x 3072    13 x 13 x 1024    13 x 13 x 125

conv    conv

**Bounding Boxes**
**+**
**IoU Scores**
**+**
**Classification Scores**

car: 0.891

**Feature**
**Alignment**
**Loss (   )**

**Object**
**Detection**
**Loss (   )**

Input: left and right spectrograms

MLP

1 x 1 x 1024

10 conv layers

3 deconv layers

13 x 13 x 125

conv

car: 0.822

Meta-data

MLP

1 x 1 x 128

**Student Network:**
**StereoSoundNet**

**Bounding Boxes**
**+**
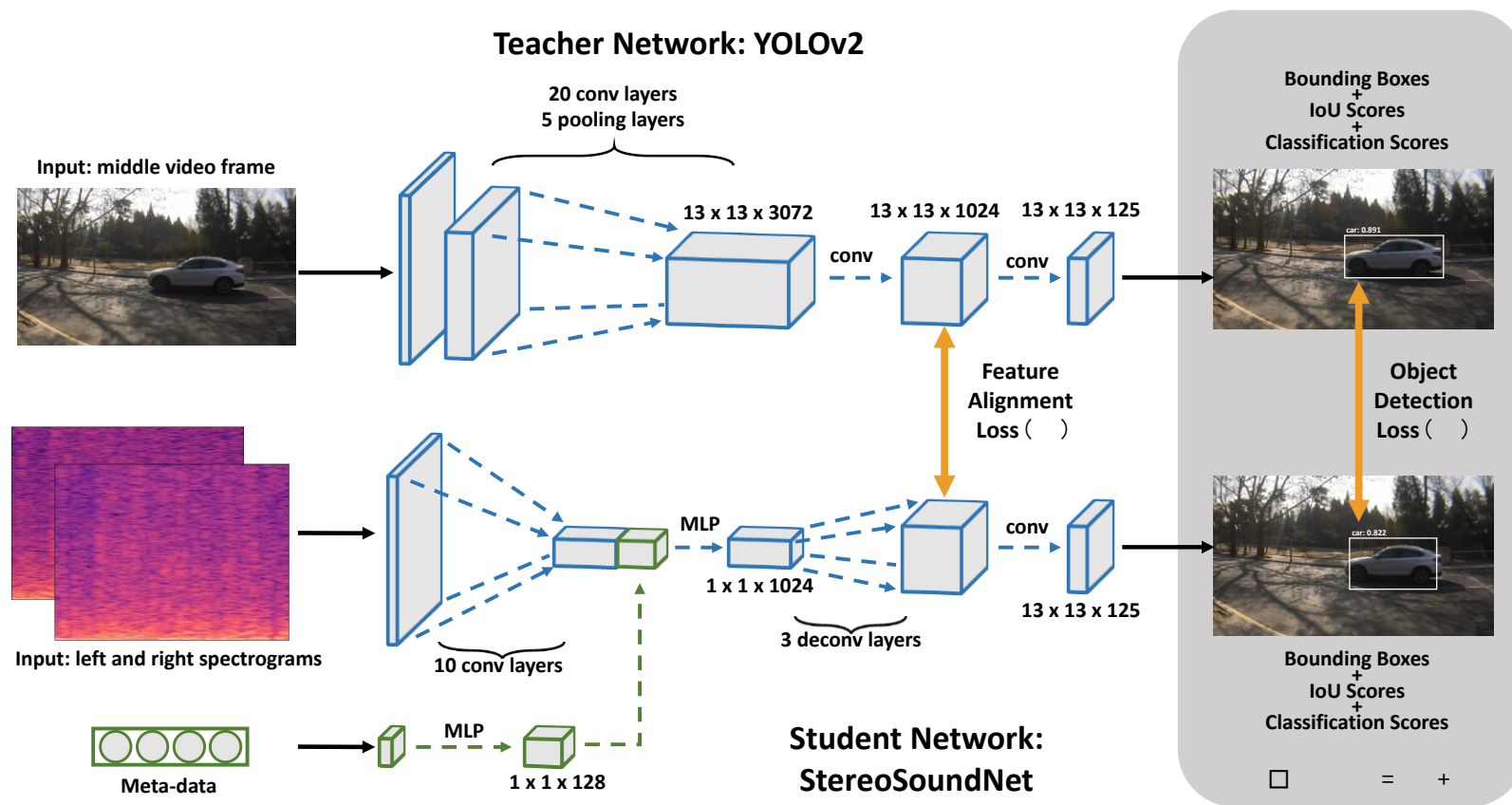**IoU Scores**
**+**
**Classification Scores**

=    +

# Our Methods

- We first get **pseudo localization labels** and **visual features** from a pre-trained YOLO

# Our Methods

- We then train a sound branch using
  - Feature alignment
  - Object alignment

# Datasets

- We have collected a dataset on diverse scenes

# Results

- Visual tracking fails in many situations:
  - Occlusion
  - Backlighting
  - Reflection on the windows
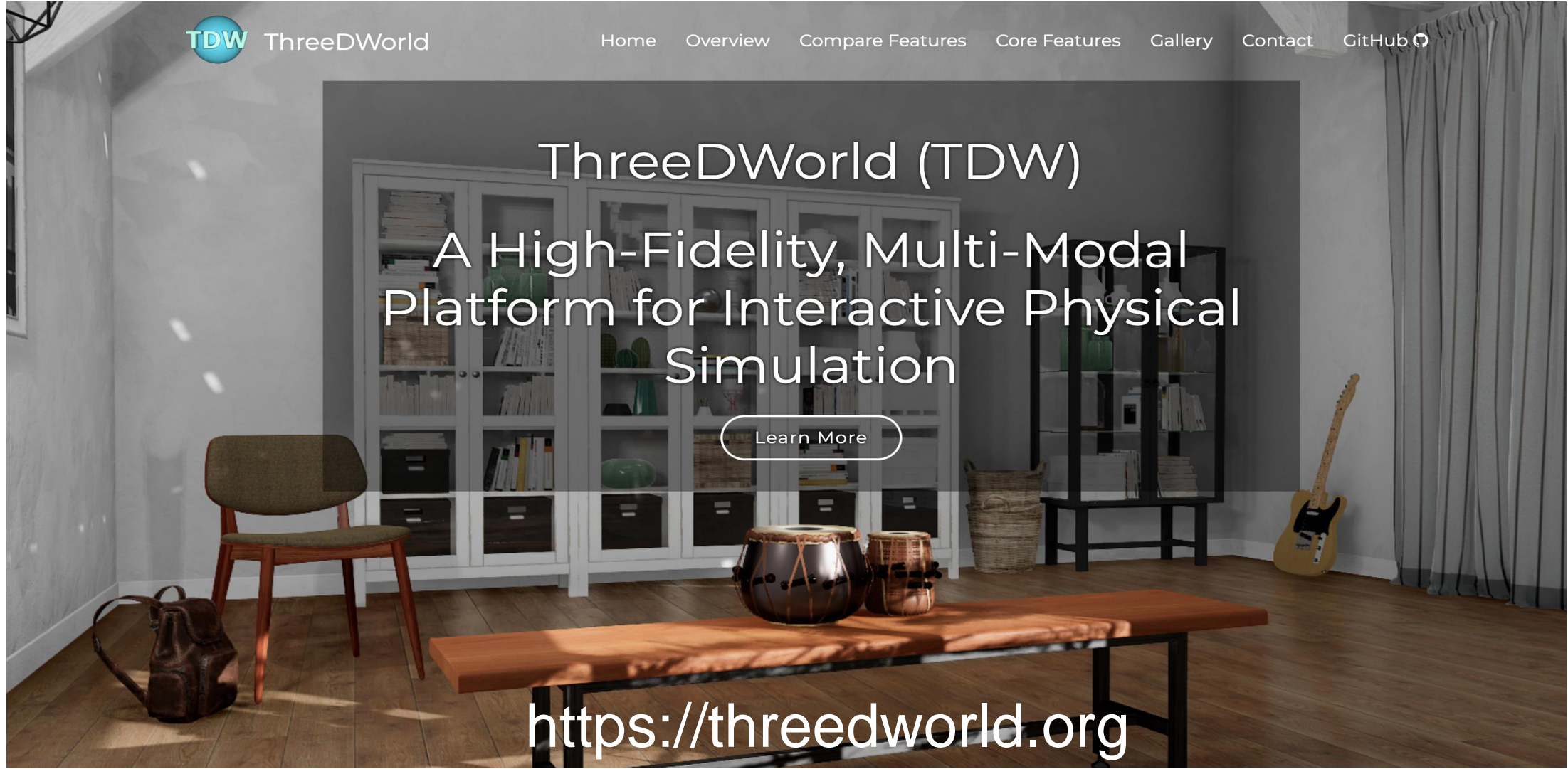  - Night scene

# Results



**Sound tracking**
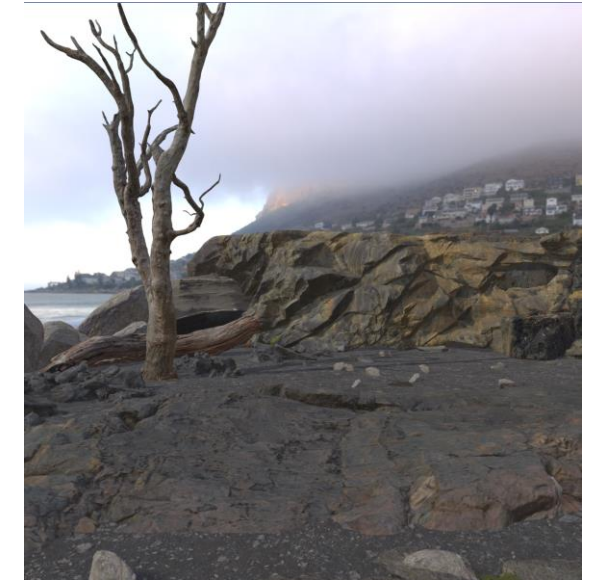
# Tracking Under Poor Lighting



**Sound Tracking**

**Visual Tracking**

# TheeDWorld: 3D Virtual World



Gan, Jeremy, et al, Threedworld: A platform for interactive multi-modal physical simulation. NeurIPS, 21

# Vision: Photo-Realistic Rendering

# Audio: Physics-Triggered Sound

# Physics Simulation



**Rigid-body**



**Soft-body**

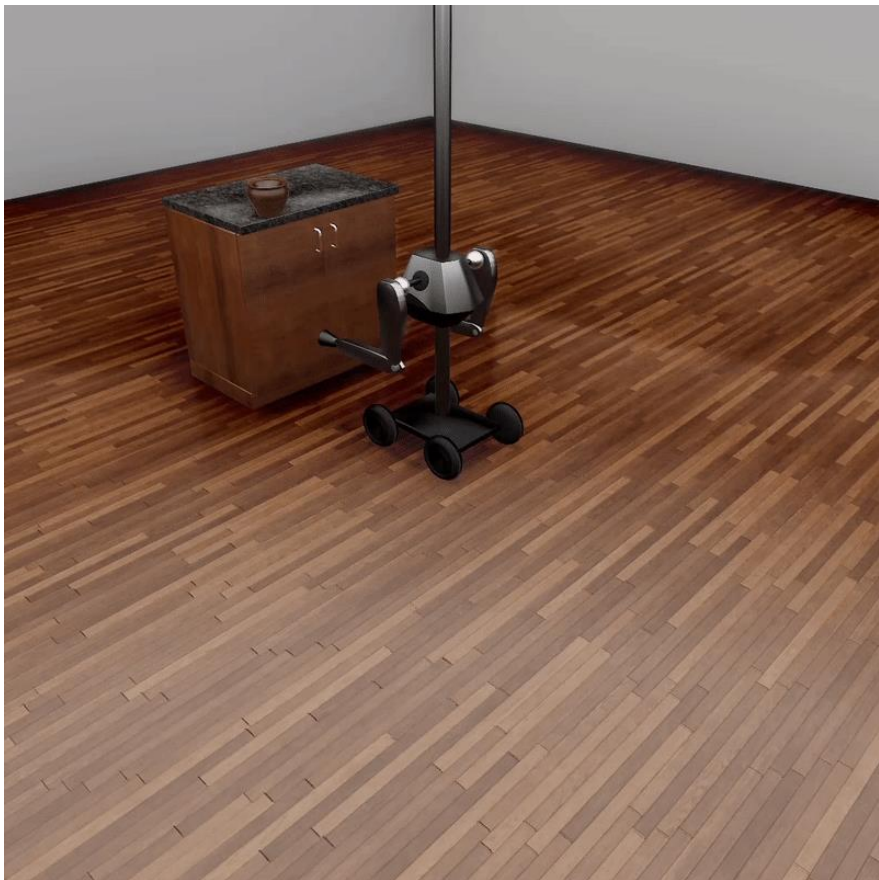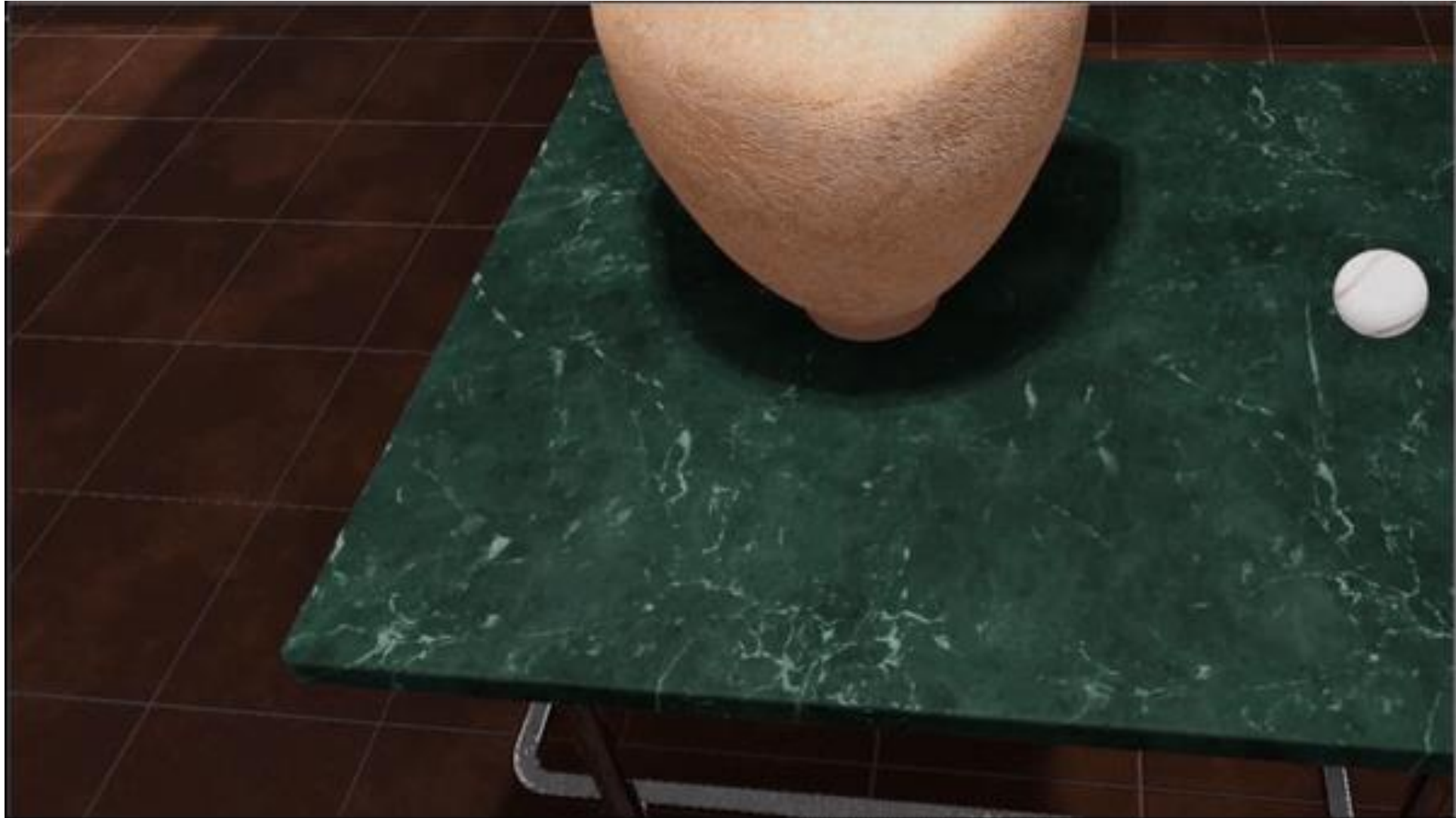# Agent Interaction



**Robot**



**Humanoid Avatar**

# Object Interaction in VR

# QA?