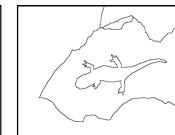
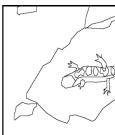
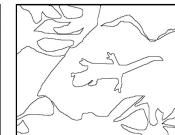


Multi-Modal AI

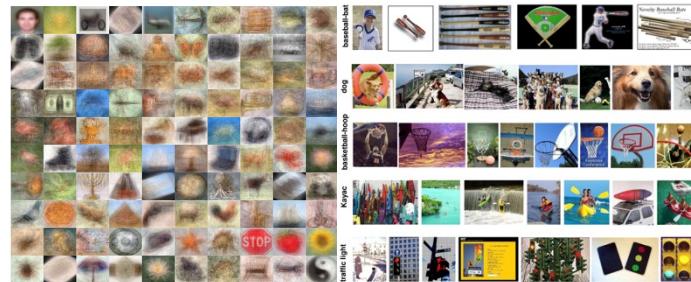
Chuang Gan



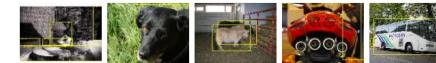
Learning from images and video frames



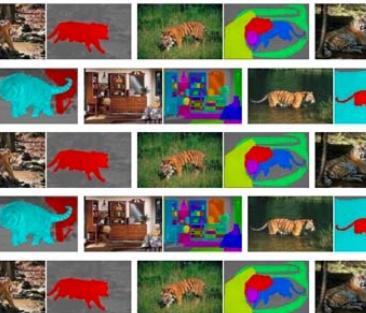
BSD (2001)



Caltech 101 (2004), Caltech 256 (2006)



PASCAL (2007-12)



LabelMe (2007)



ImageNet (2009)



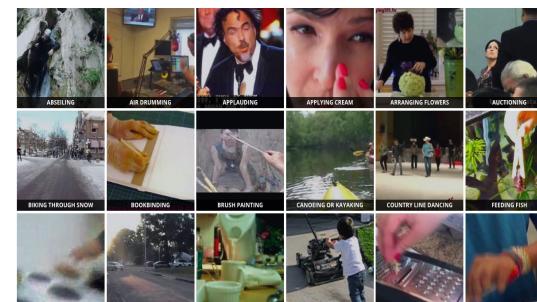
SUN (2010)



UCF-101 (2012)



Youtube-8M (2017)



Kinetics (2017)

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

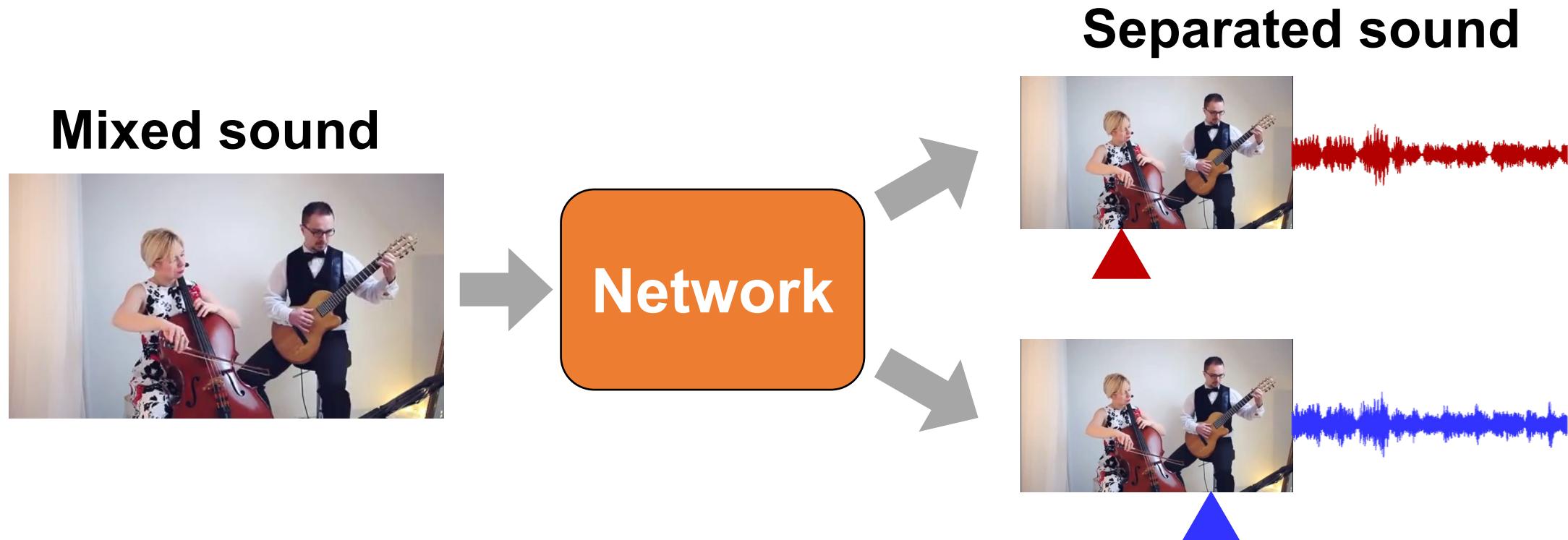
Given a music performance video...

Mixed sound



Task: Visual Sound Separation

...we aim to separate two sounds played by different instruments.

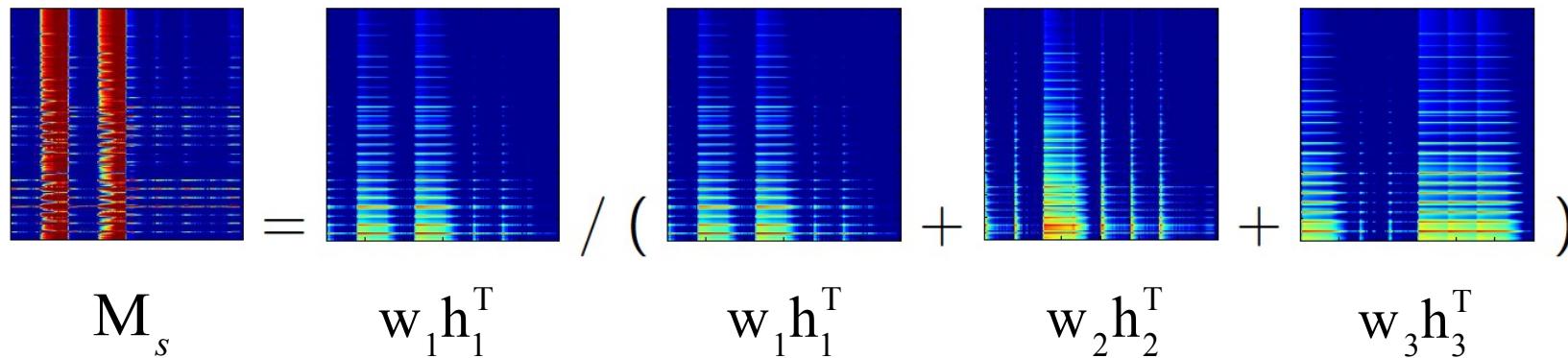


Source Separation: Traditional Approach

□ To separate one component out of K:

- 1. Estimate the mask for the target component

$$\mathbf{M}_s = \frac{\mathbf{w}_1 \mathbf{h}_1^T}{\sum_{i=1}^K \mathbf{w}_i \mathbf{h}_i^T}$$

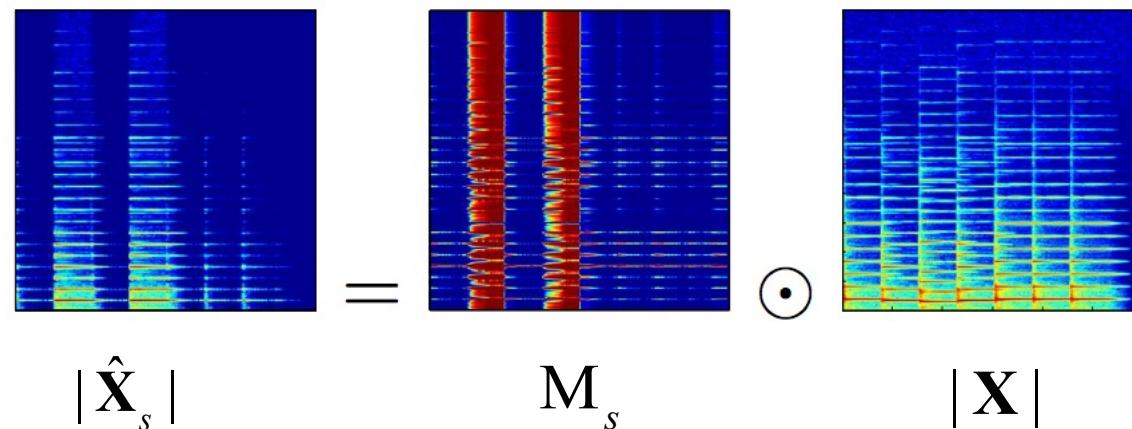


Source Separation: Traditional Approach

□ To separate one component out of K:

- 2. Masking on the input spectrogram to separate the component

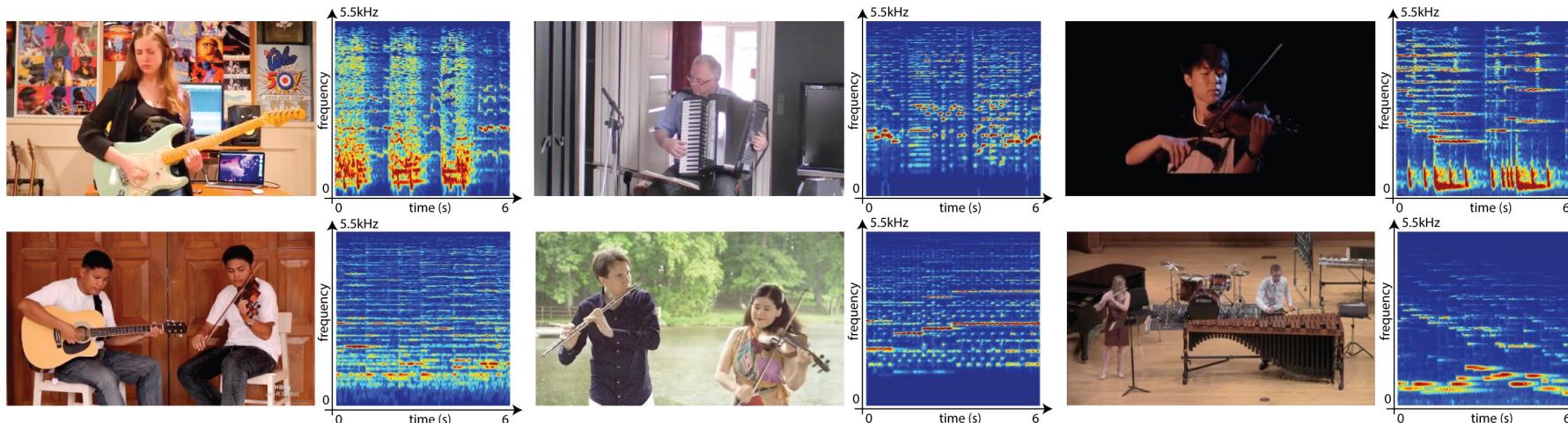
$$|\hat{\mathbf{X}}_s| = \mathbf{M}_s \odot |\mathbf{X}|$$



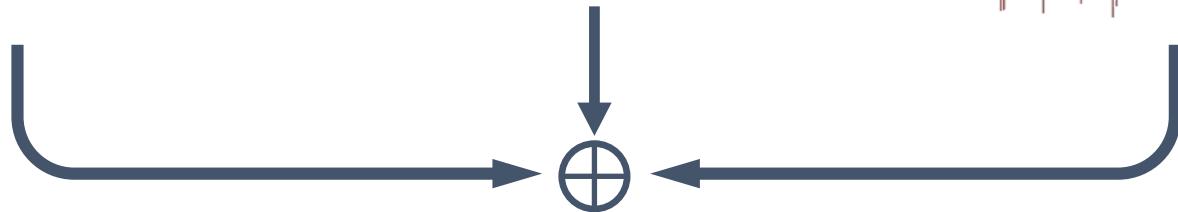
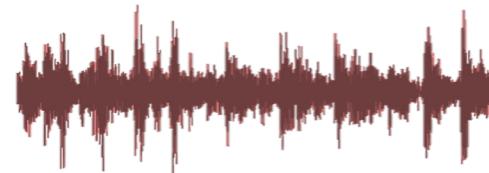
Our Ideas: Learning from Music Videos

□ Internet music videos

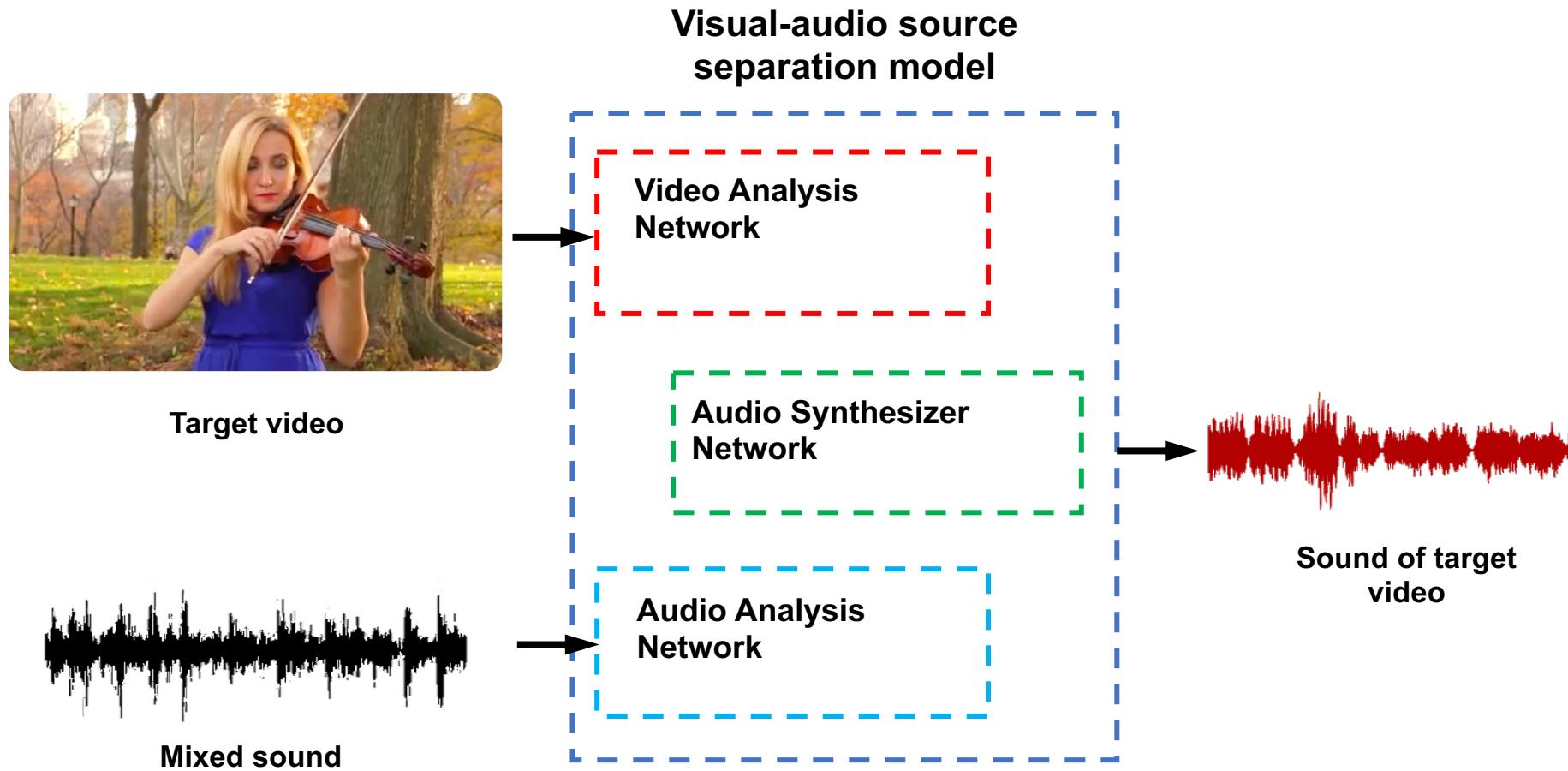
- Keyword search without labeling
- >20 kinds of commonly seen musical instruments
- >1000 solos and duets



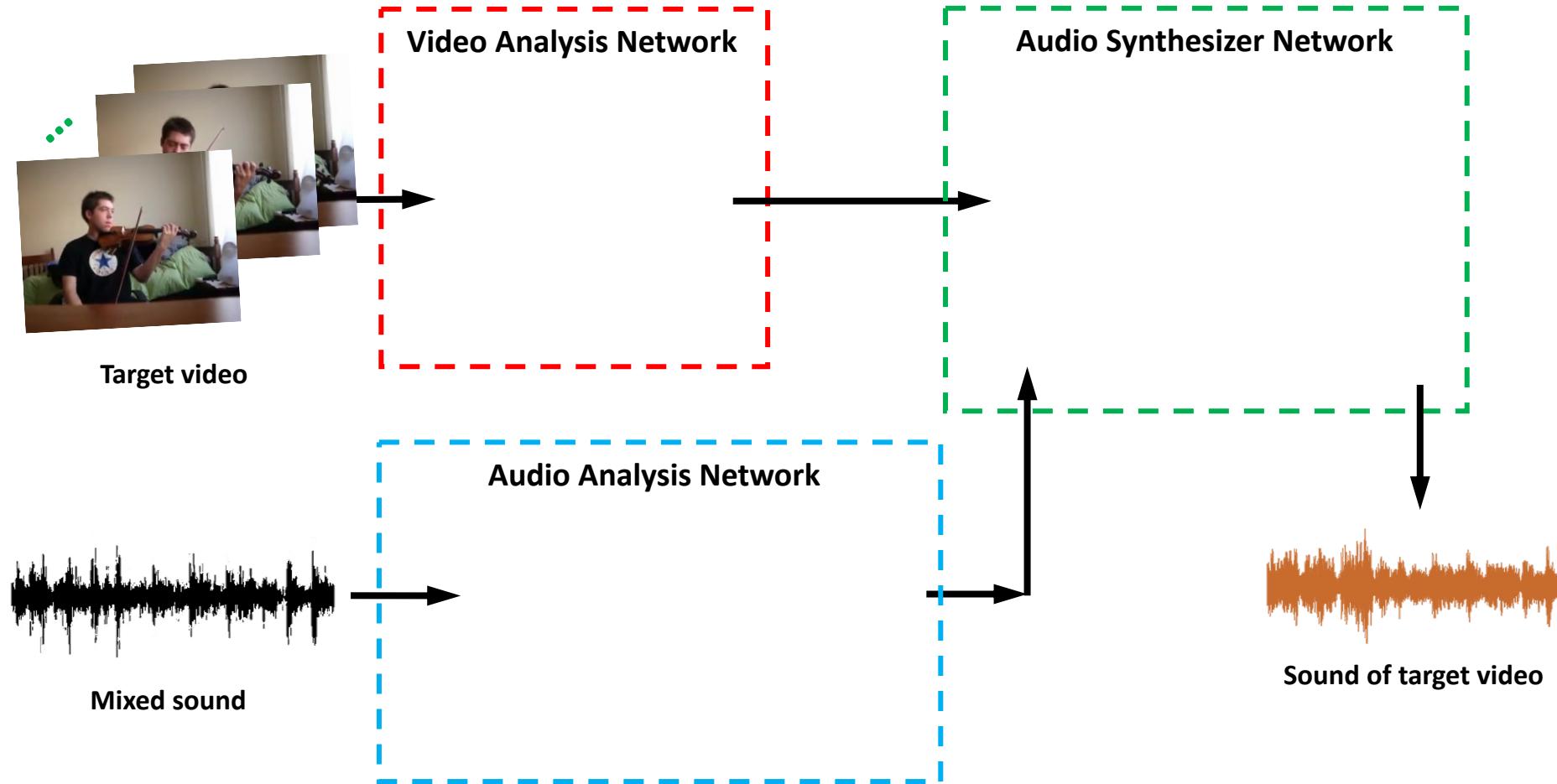
Mixing the Sound



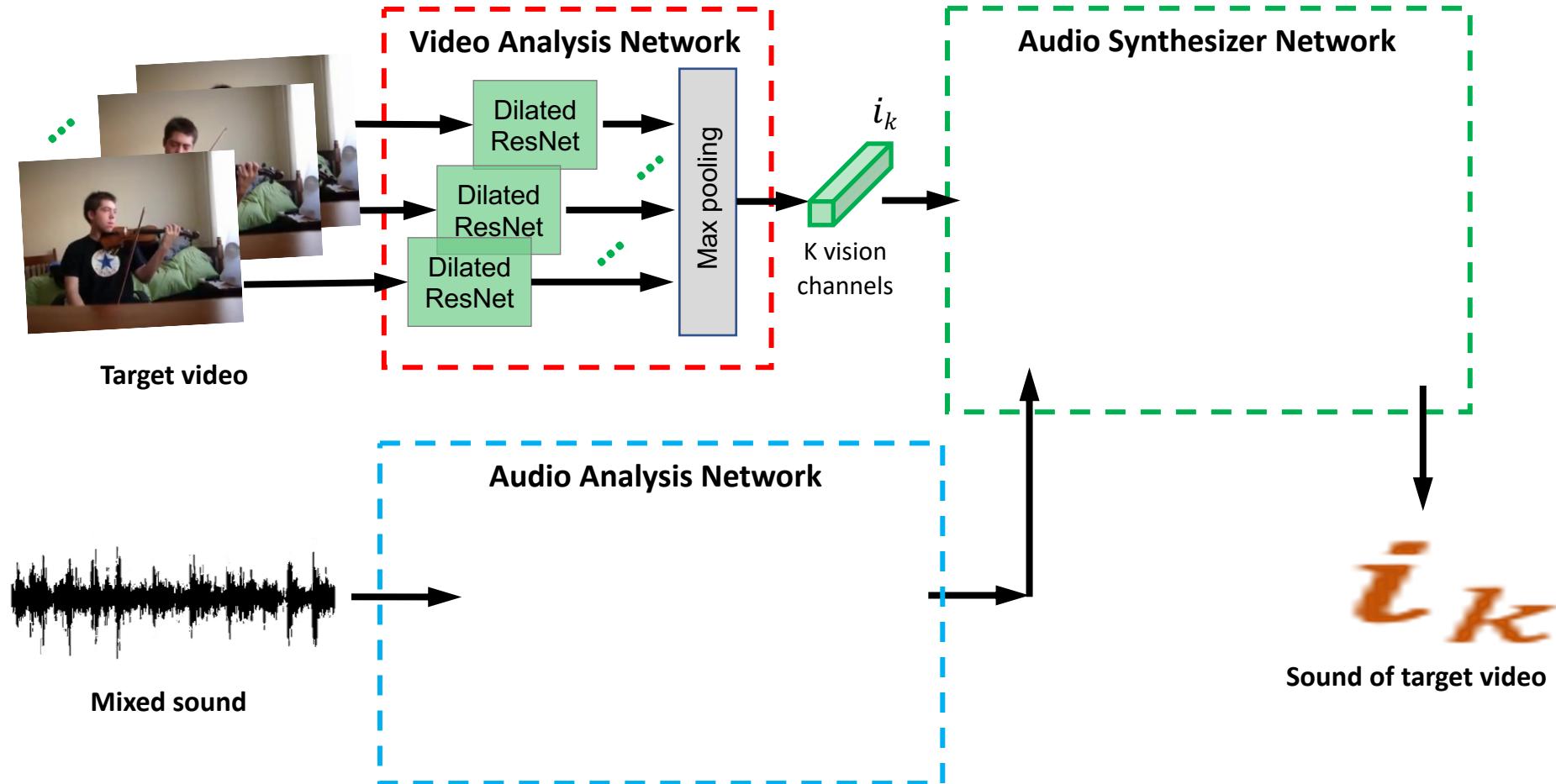
Vision to Rescue for Self-supervised Prediction



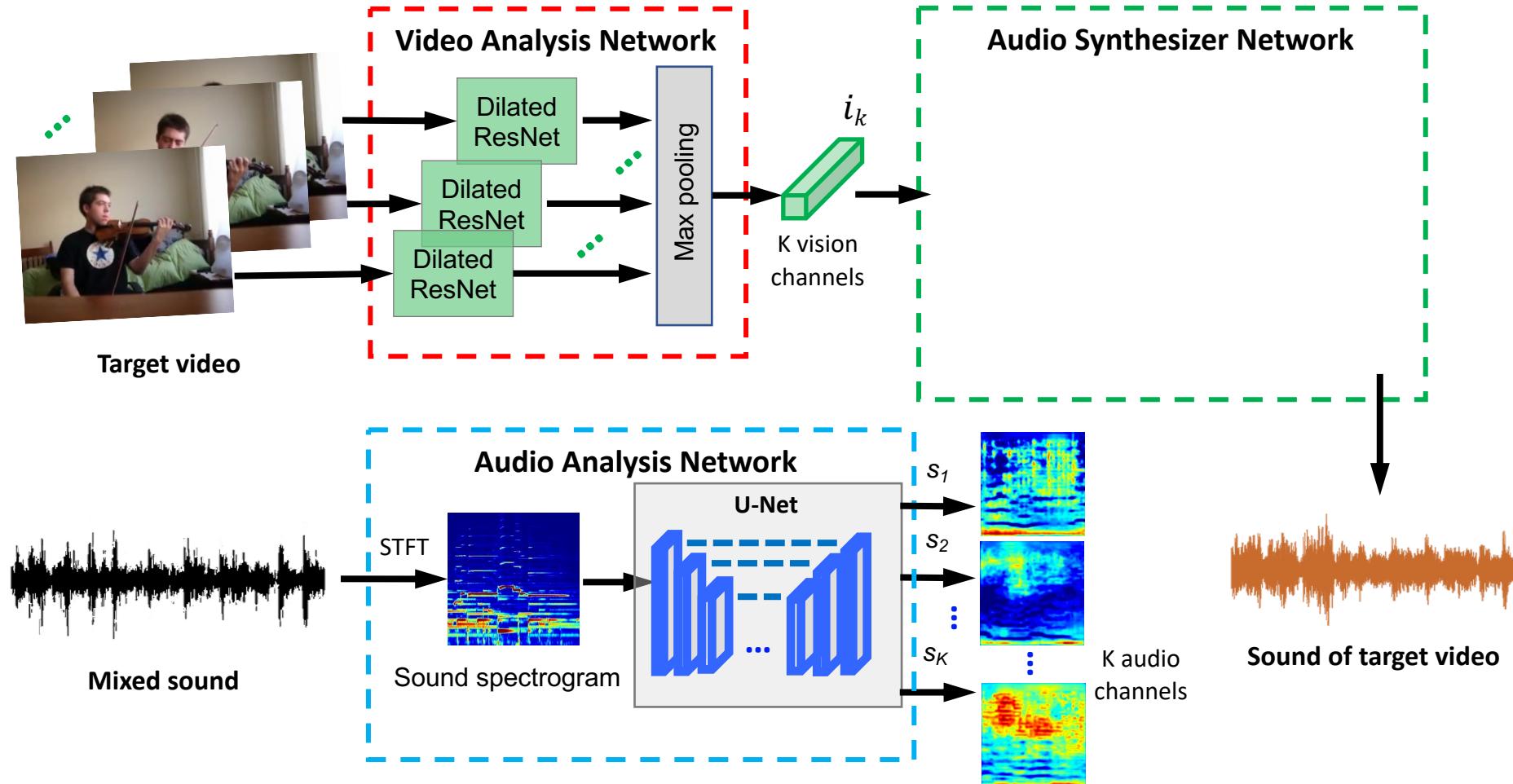
Mix-and-Separate Framework



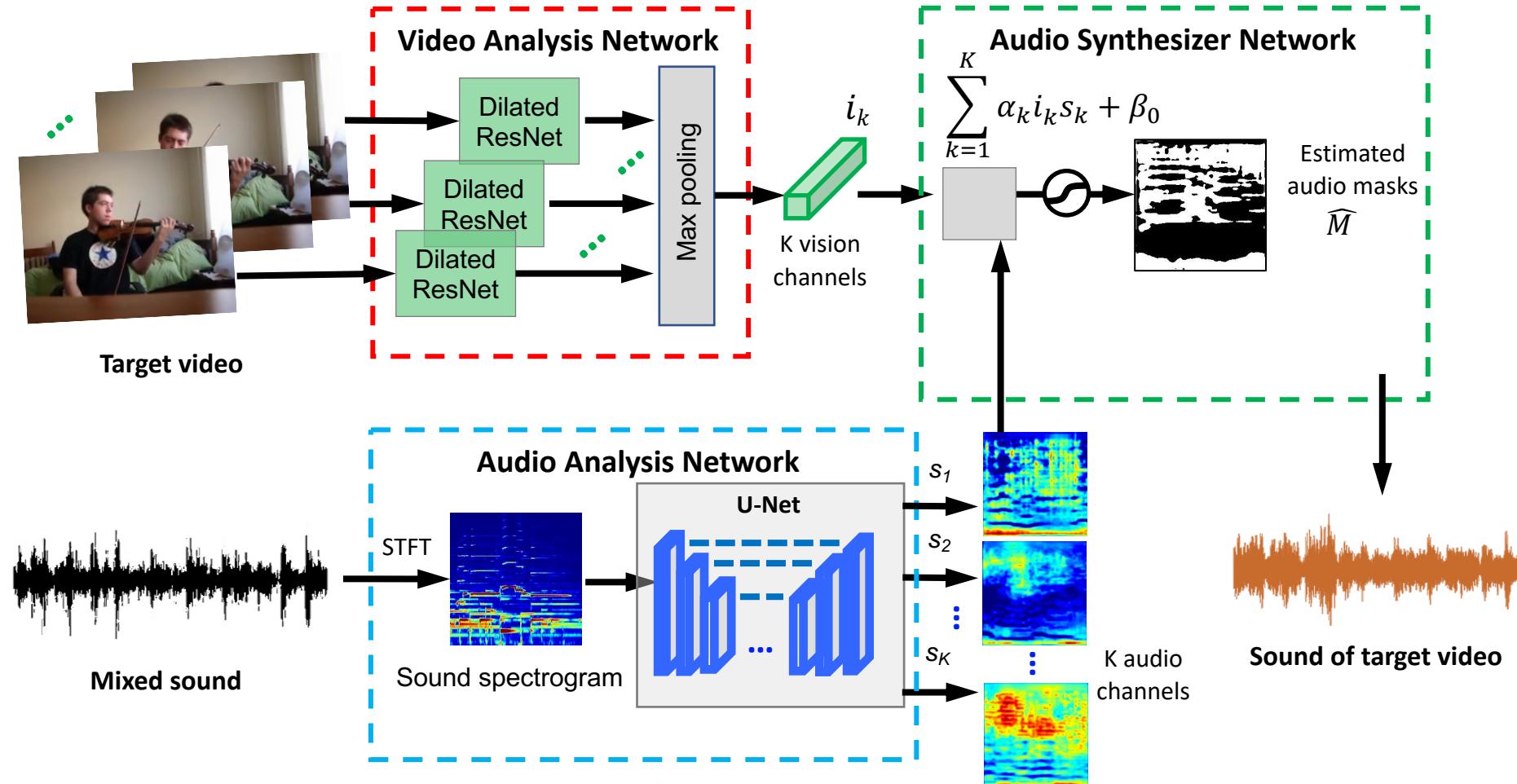
Mix-and-Separate Framework



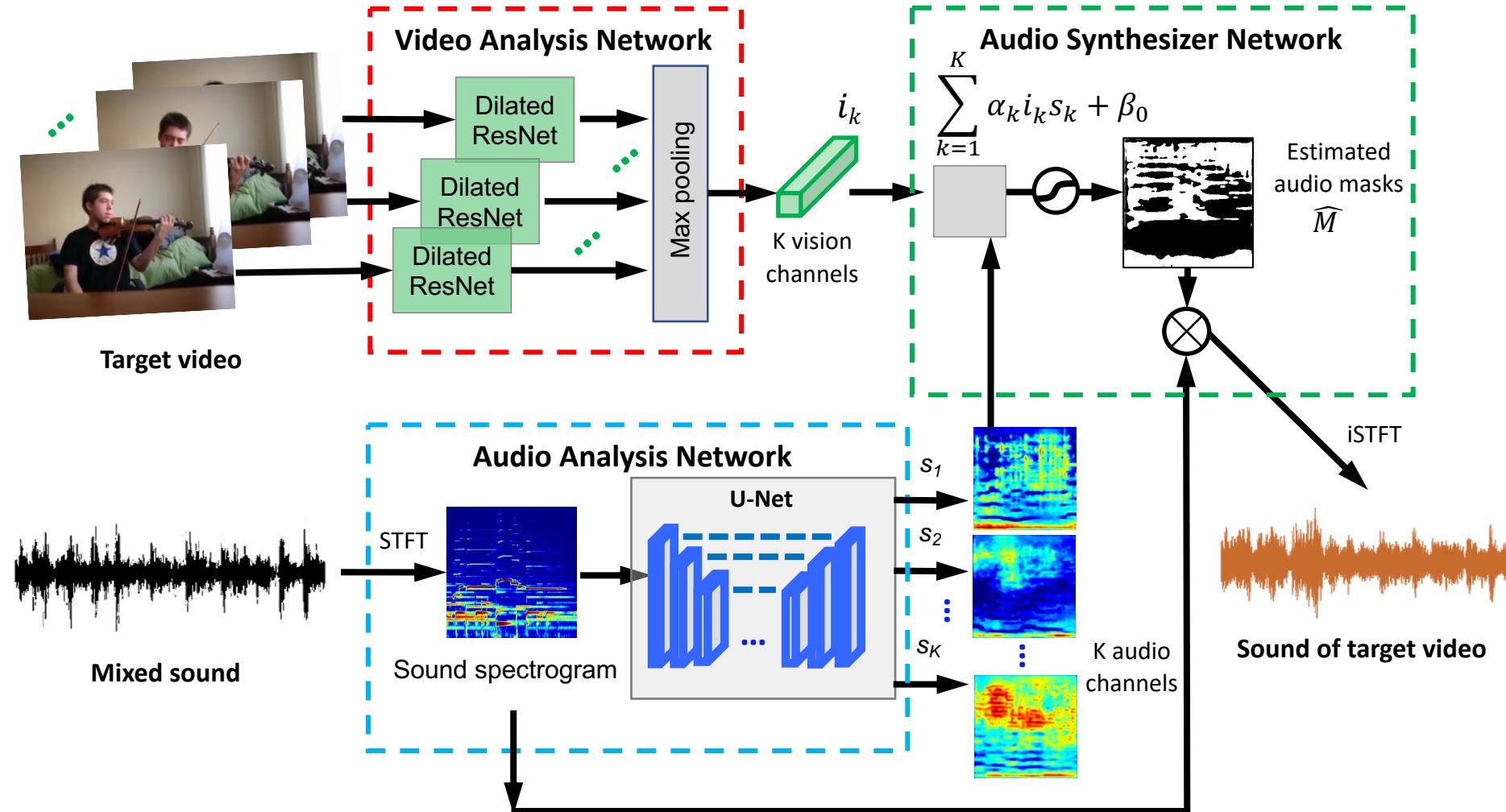
Mix-and-Separate Framework



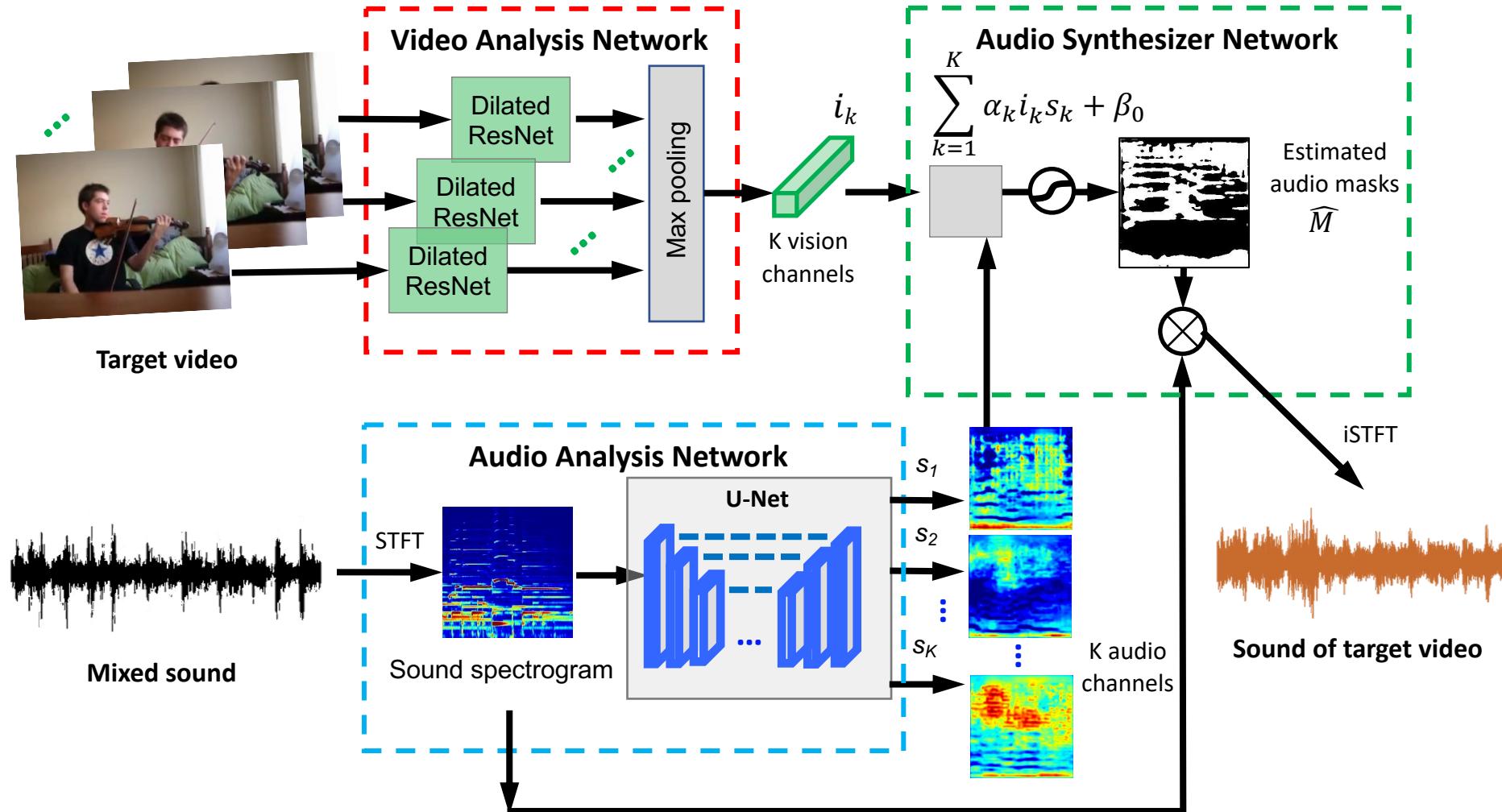
Mix-and-Separate Framework



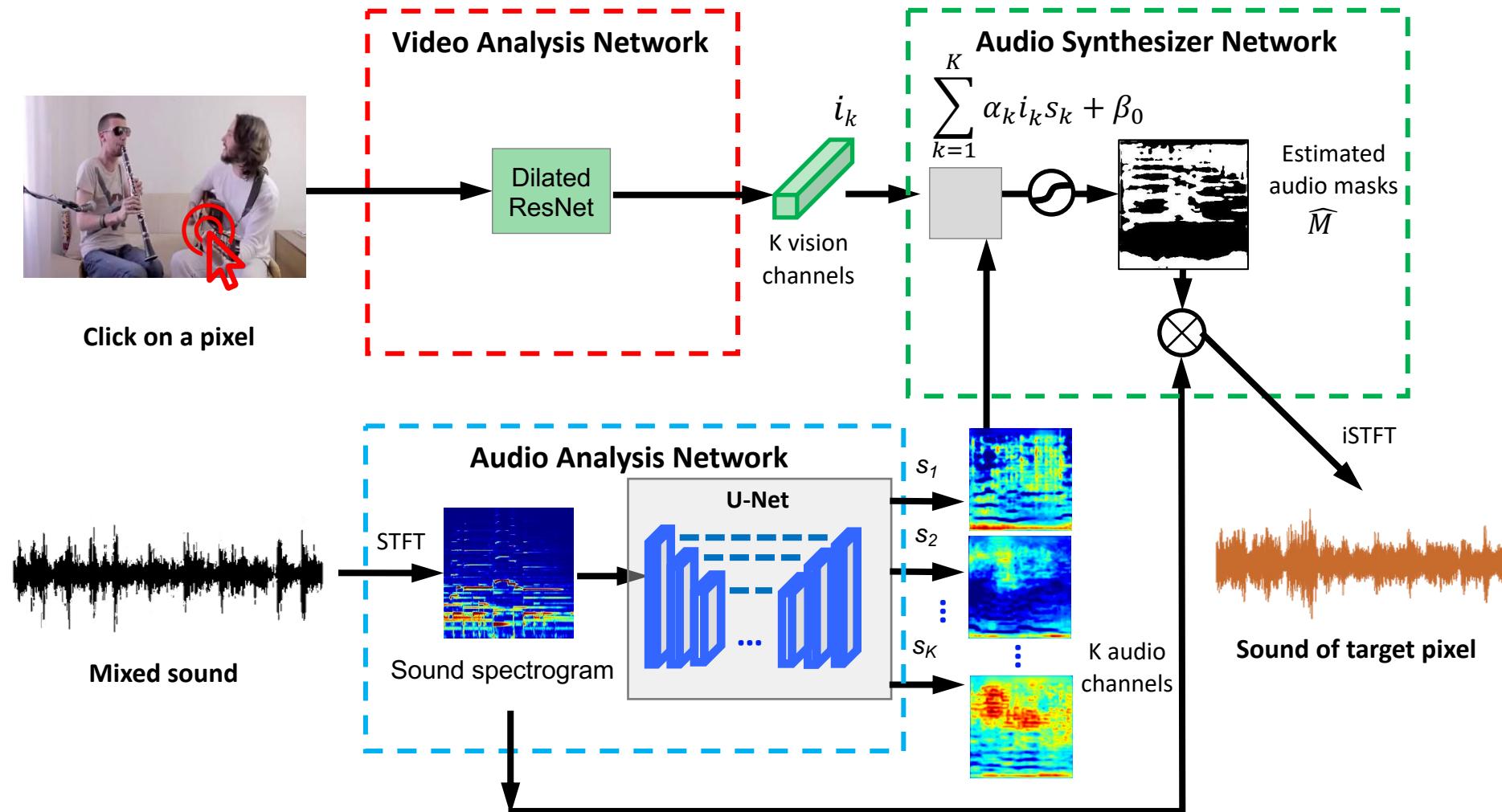
Mix-and-Separate Framework



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...



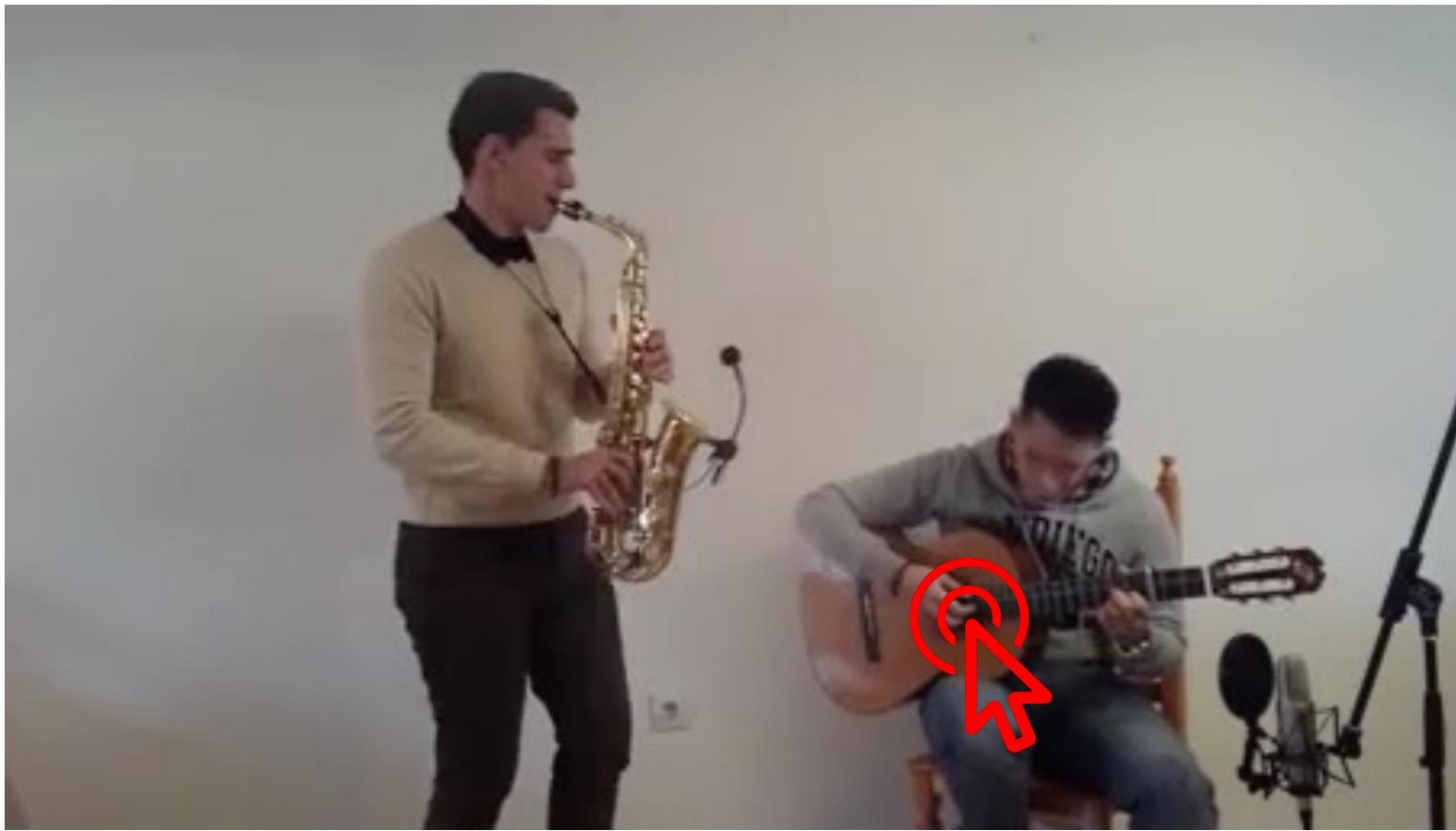
Original Video



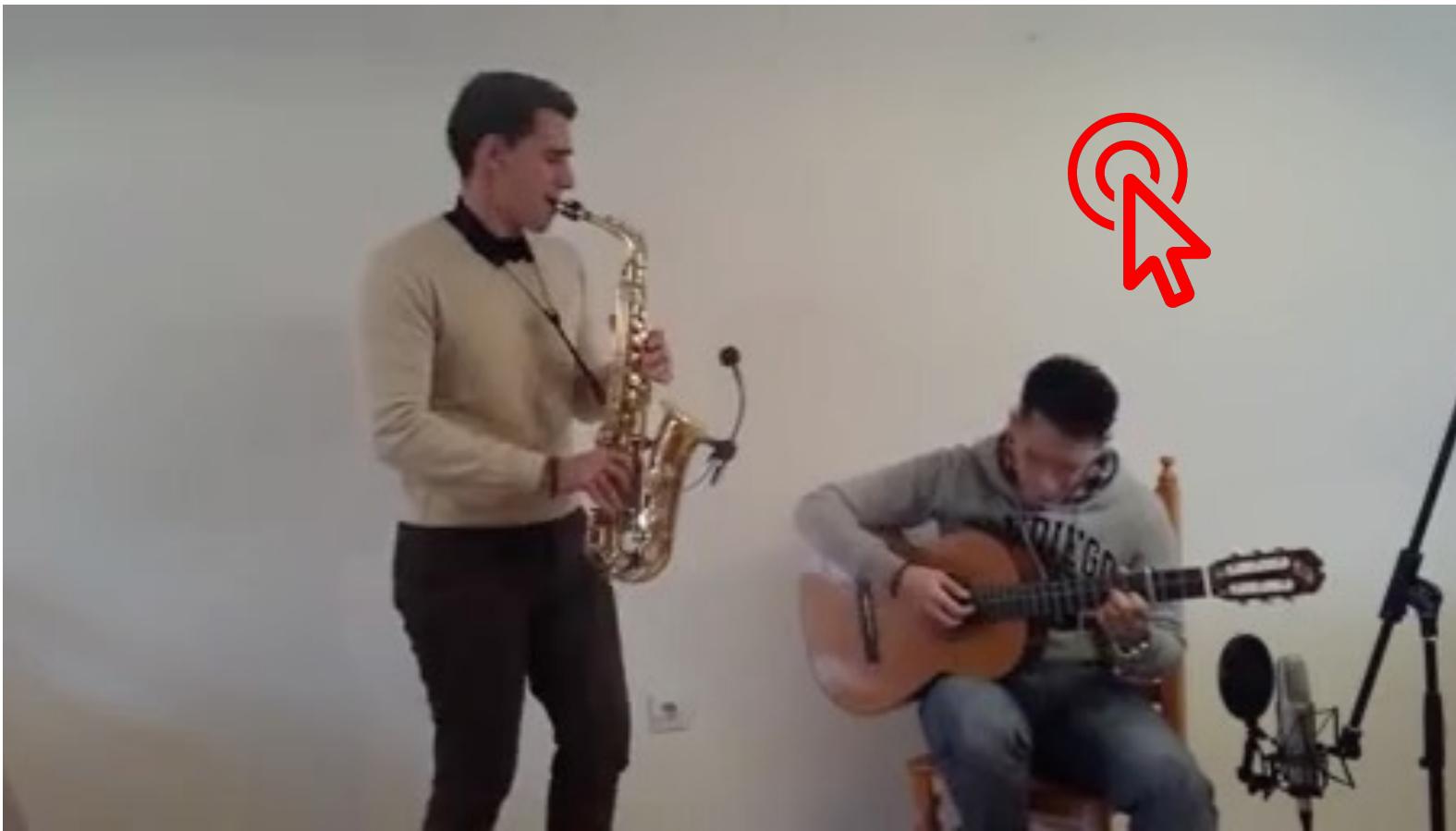
The sound of clicked object...



The sound of clicked object...



The sound of clicked object...



Application: Music Remix



Application: Music Remix



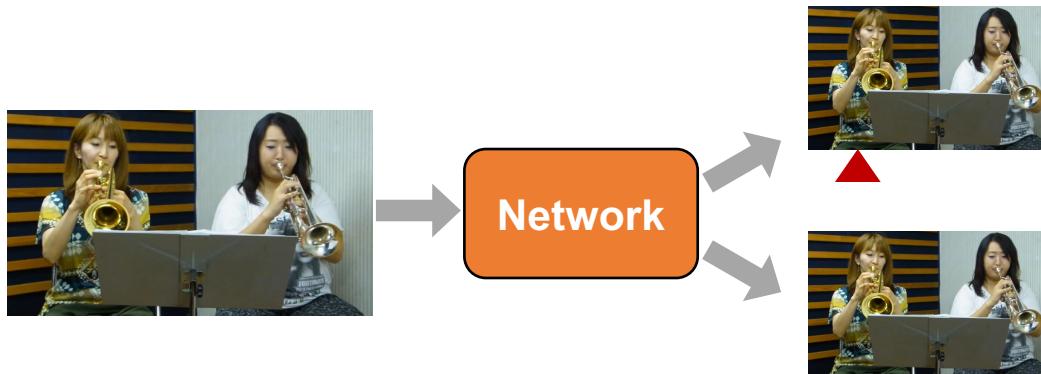
Play



Limitation

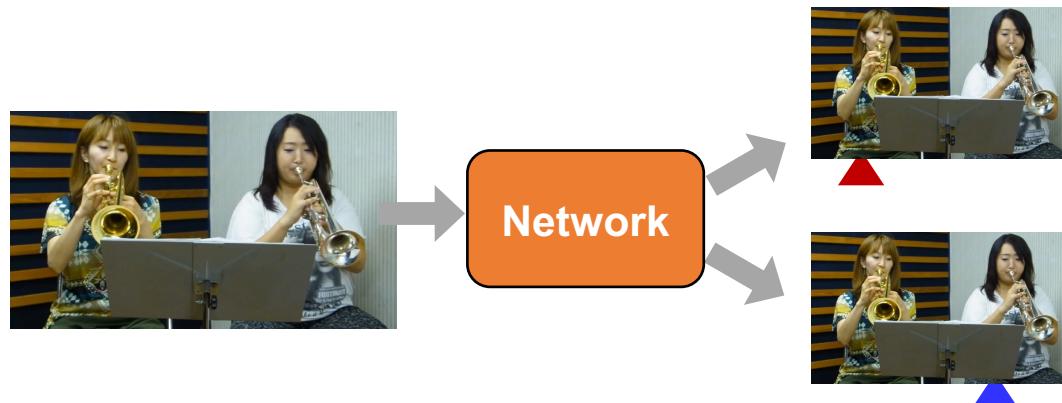
- Most existing methods use **raw pixel** or **optical flow** as input.

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



Our Ideas

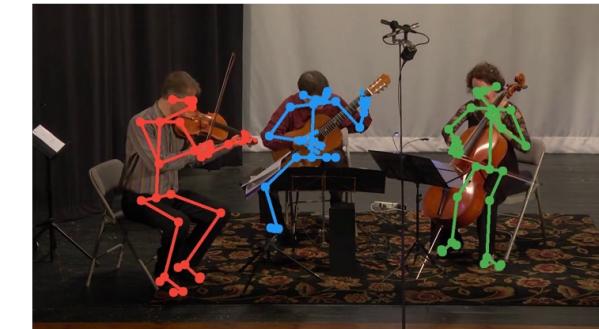
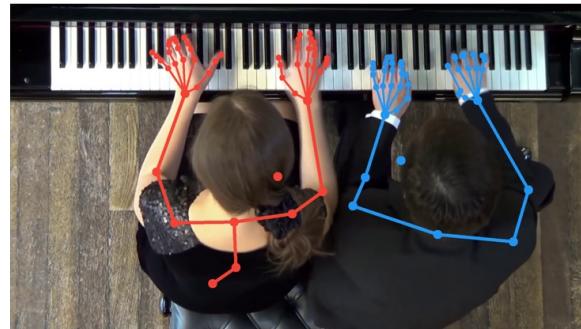
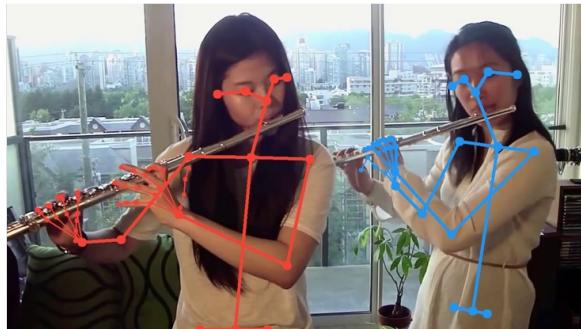
- Problem: limited to **separate multiple instruments of the same types.**



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

Music Gesture

❑ Keypoint-based structured representations



Visual sound separation results

Sound of Motion



Mixed sound



Separated
sound1



Separated
sound2

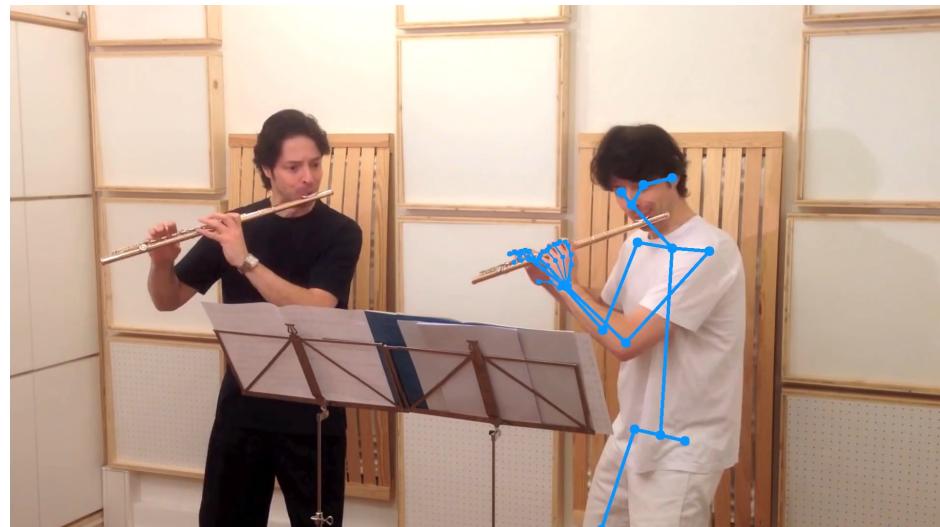
Music Gesture



Mixed sound



Separated
sound1



Separated
sound2

Sound of Motion



Mixed sound



Separated
sound1



Separated
sound2

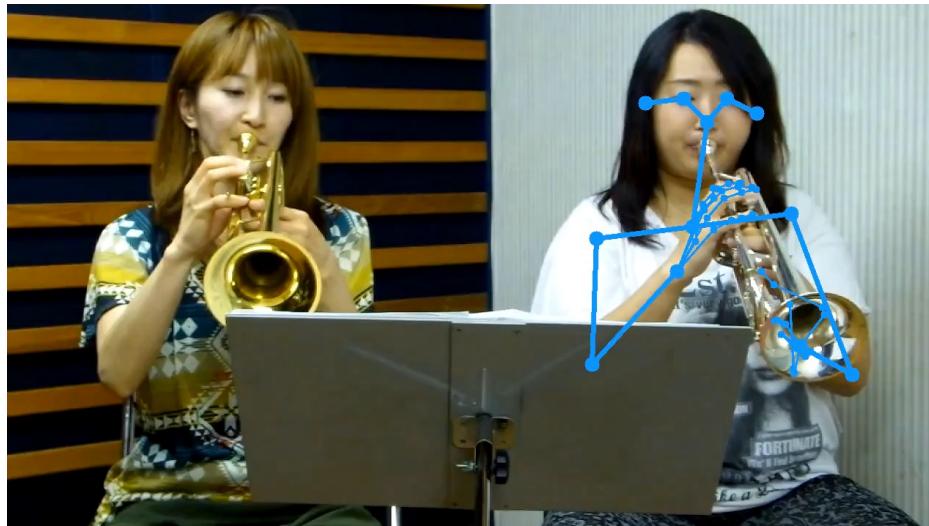
Music Gesture



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

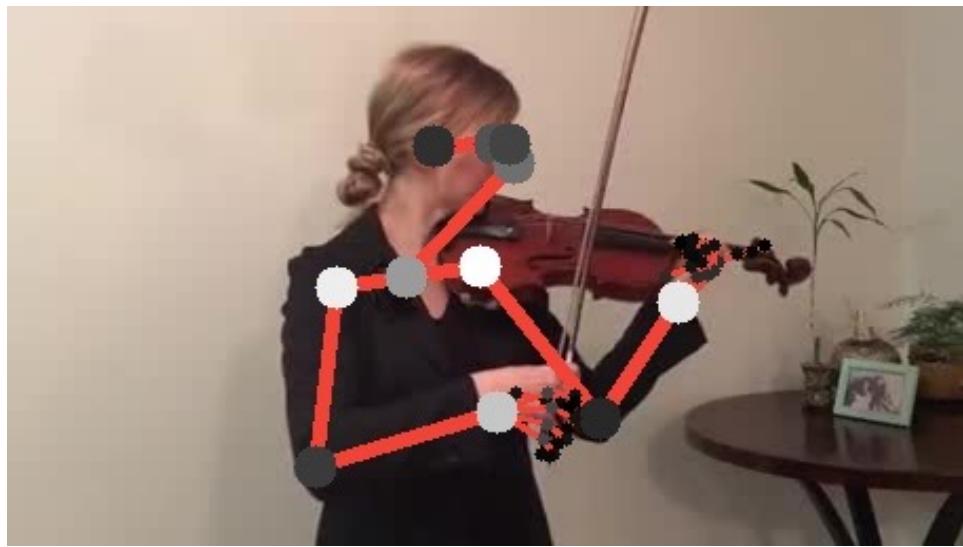


Separated sound3



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.



Silent music performance video

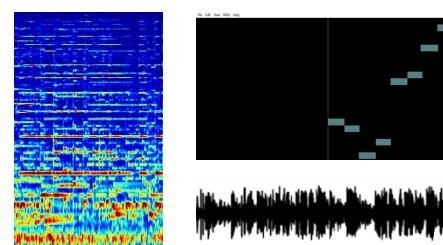
Deep Neural Network



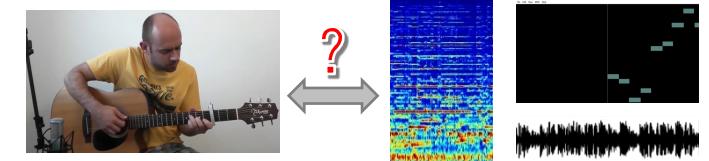
Performance with generated sound

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 ?



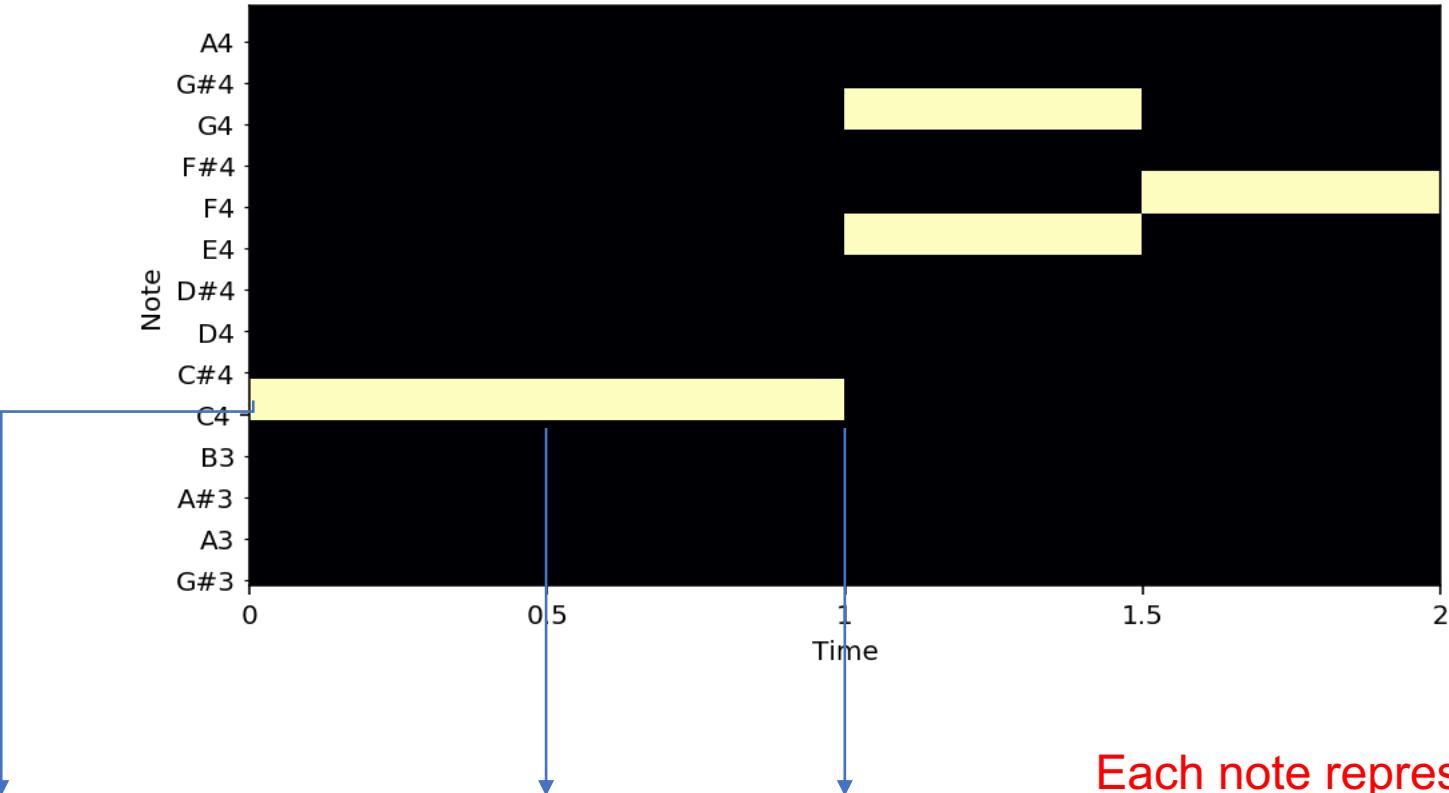
Challenges

- 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



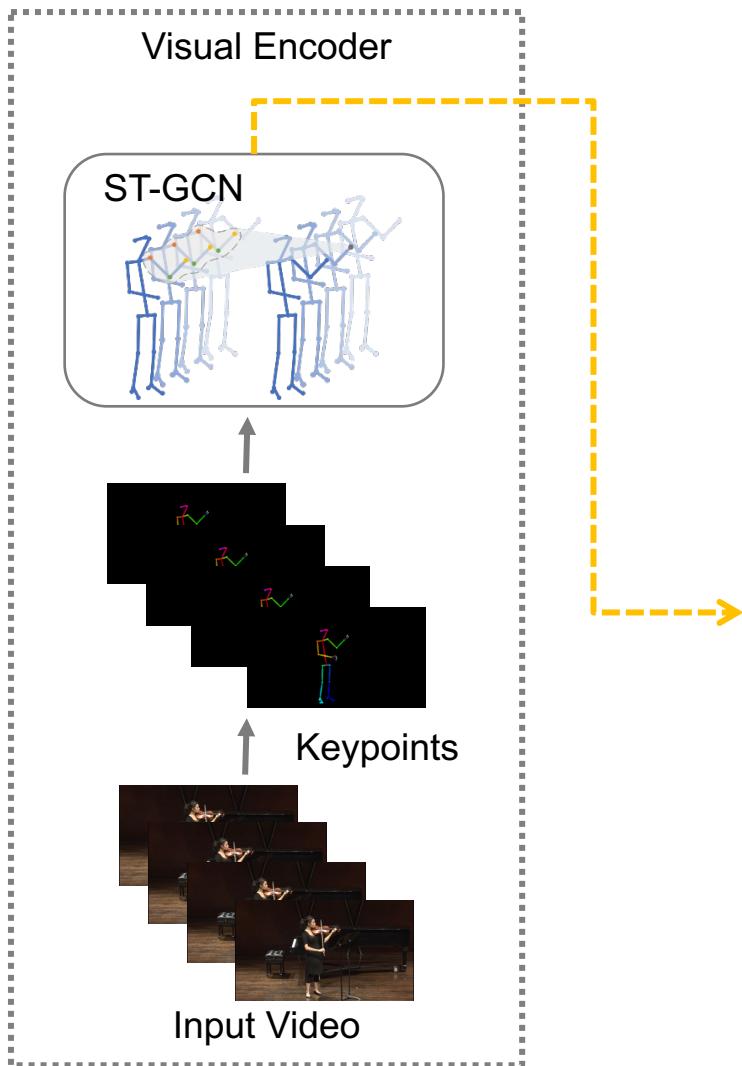
Velocity Event, Note On Event Time Shift Event Note Off Event
201 39 237 127

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

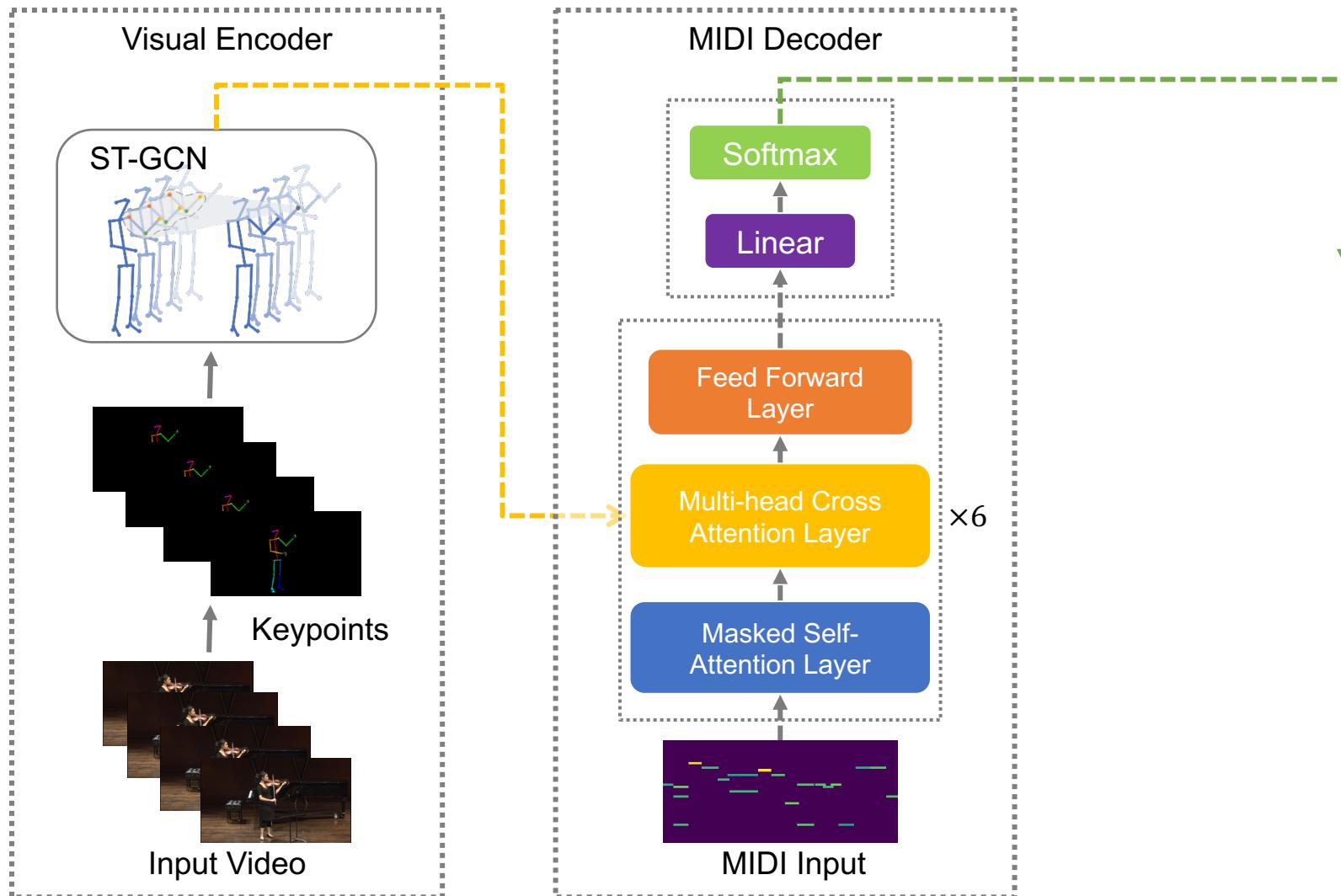
Our method

We first encode **human keypoints** by using graph CNN.



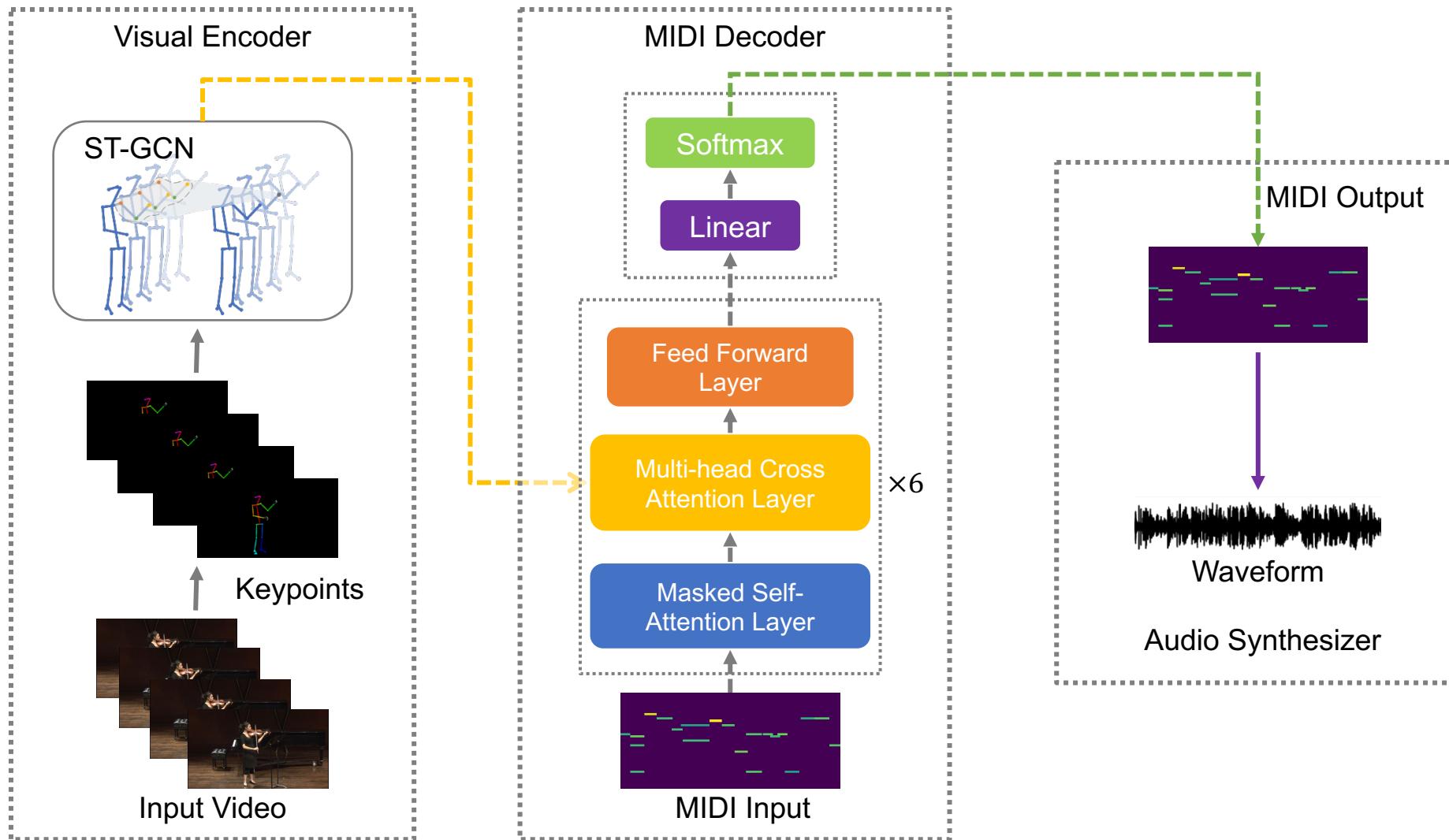
Our method

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



Original prediction

Style editing



A major



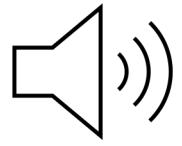
F major



G major

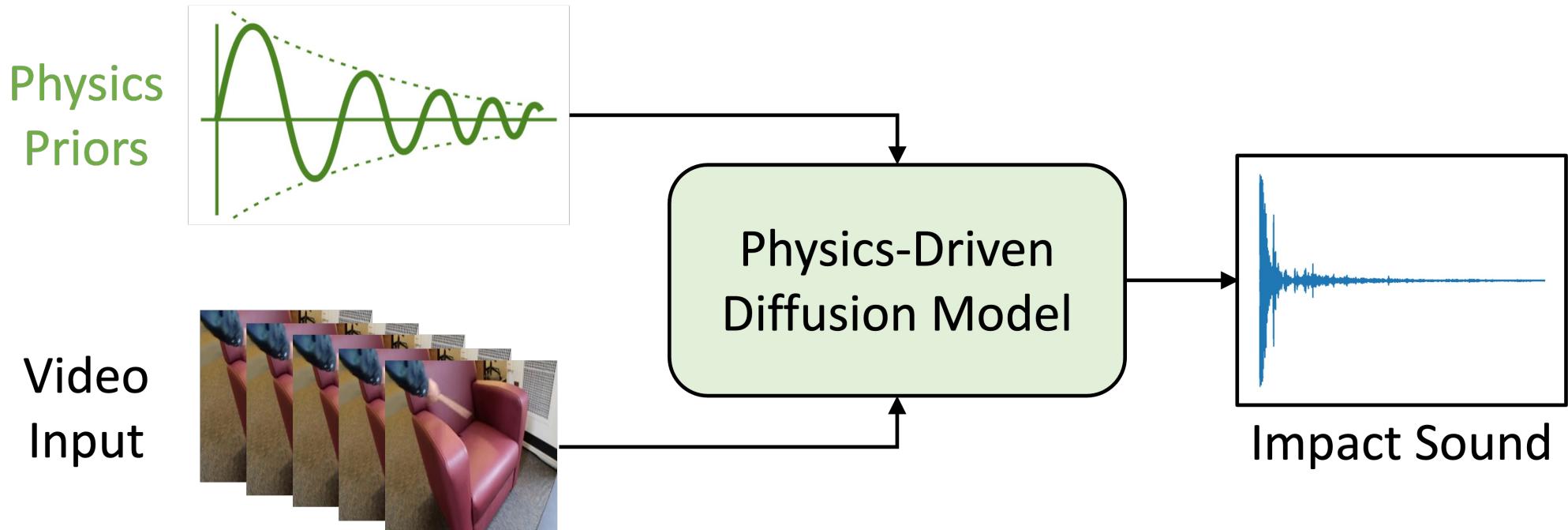


An impact sound of
physical object interactions
is critical to perception



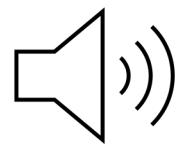
Could we generate the
impact sound from vision?

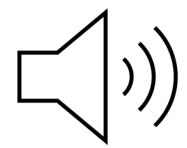
Our Framework

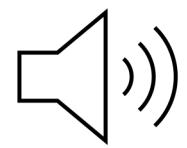
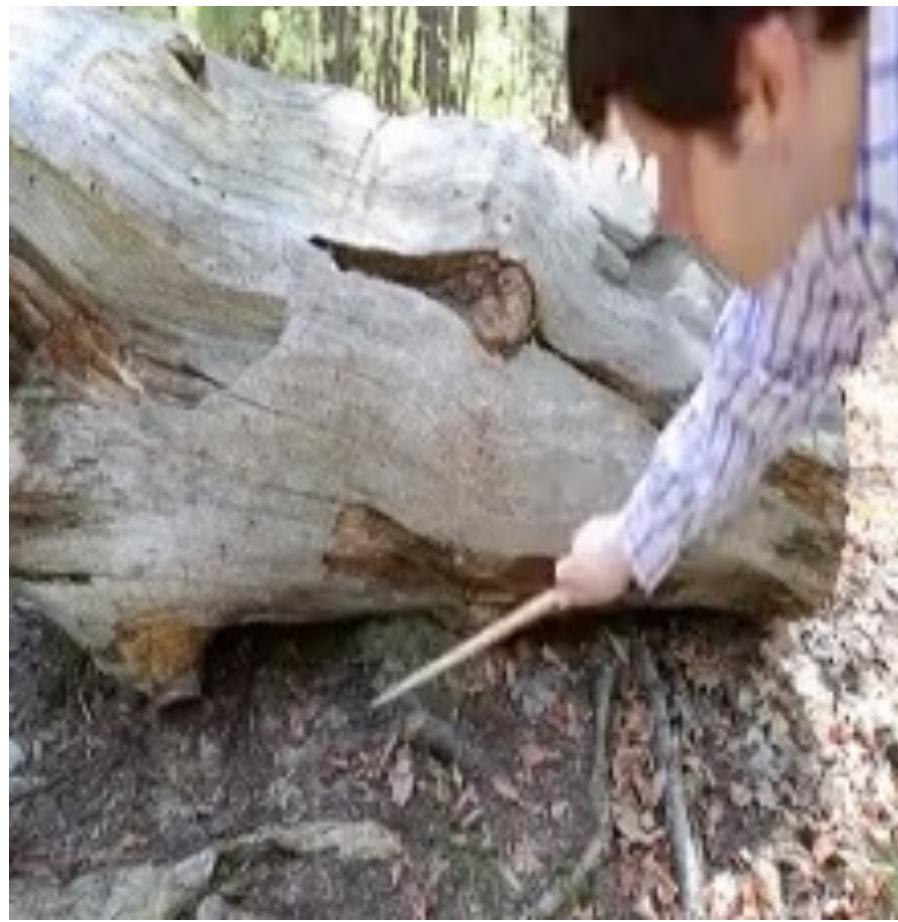


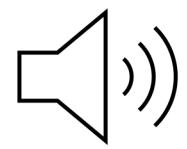


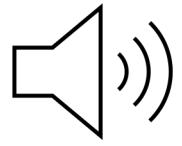
Our model is applicable to
a variety of videos and materials

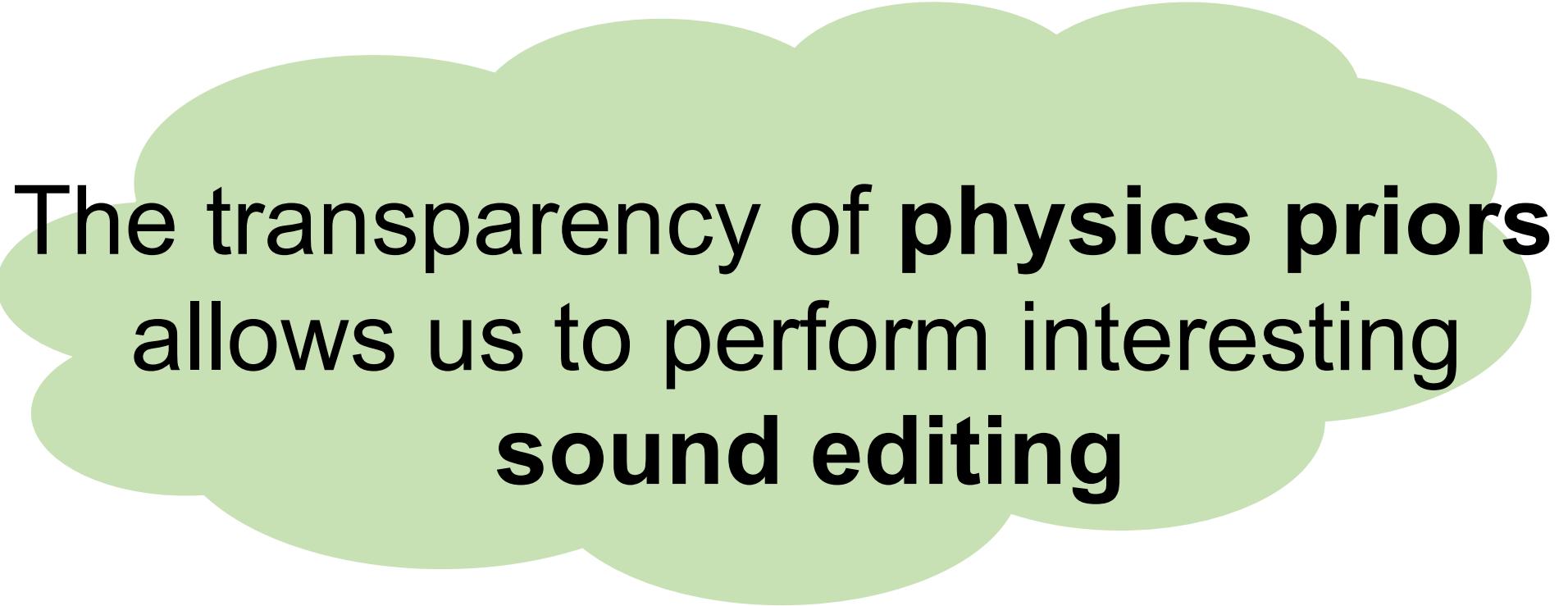






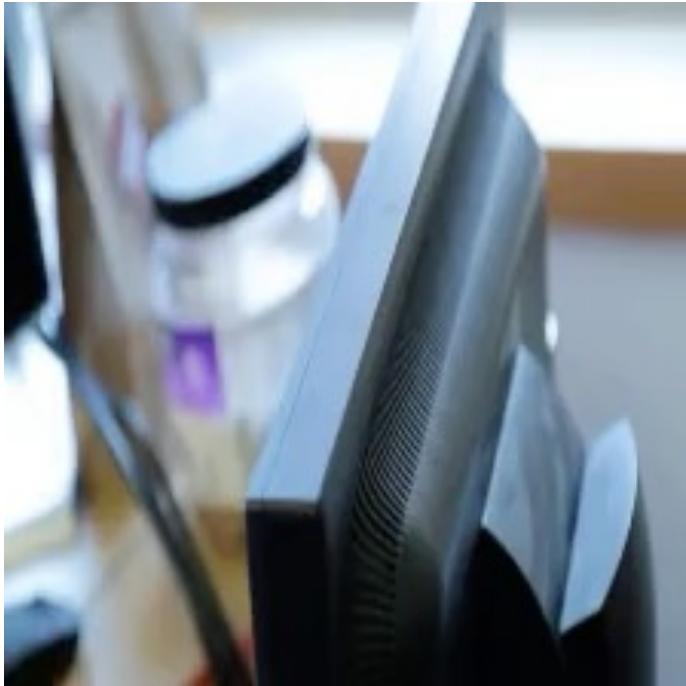






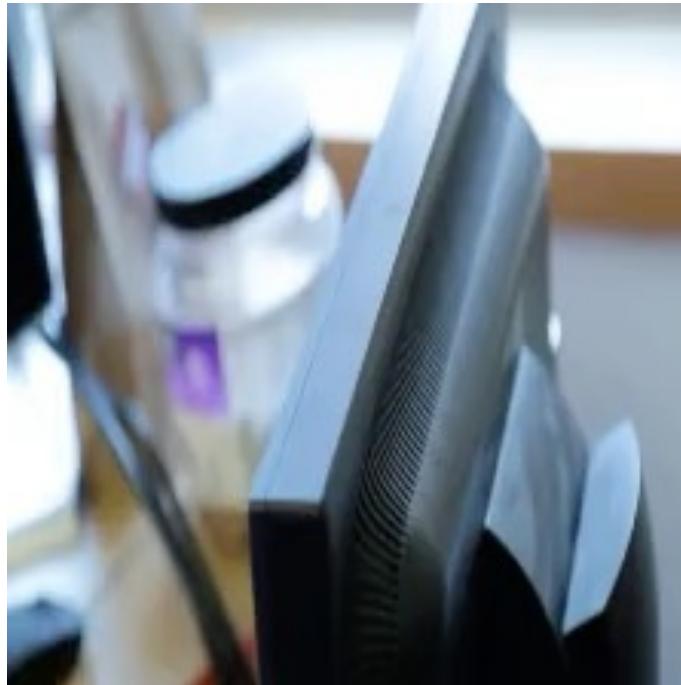
The transparency of physics priors
allows us to perform interesting
sound editing

Physics Priors of Glass + Video Input

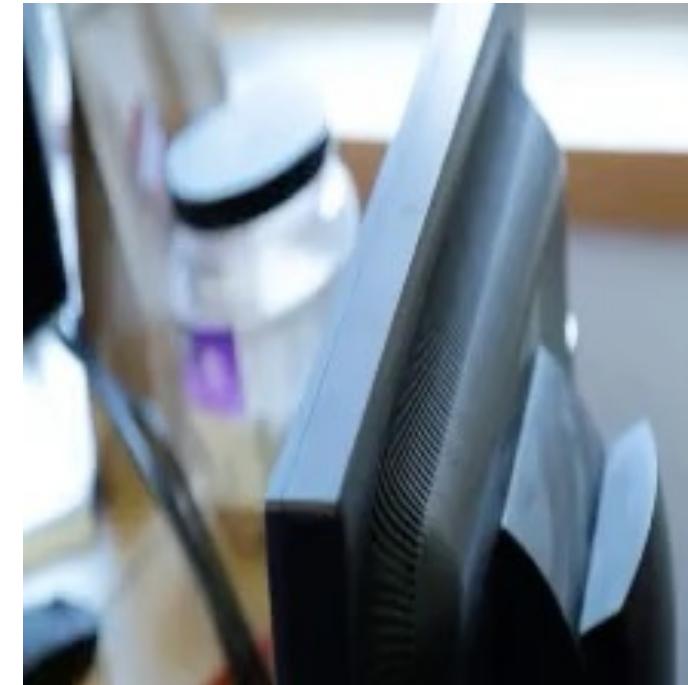


Original Result

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



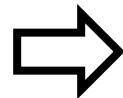
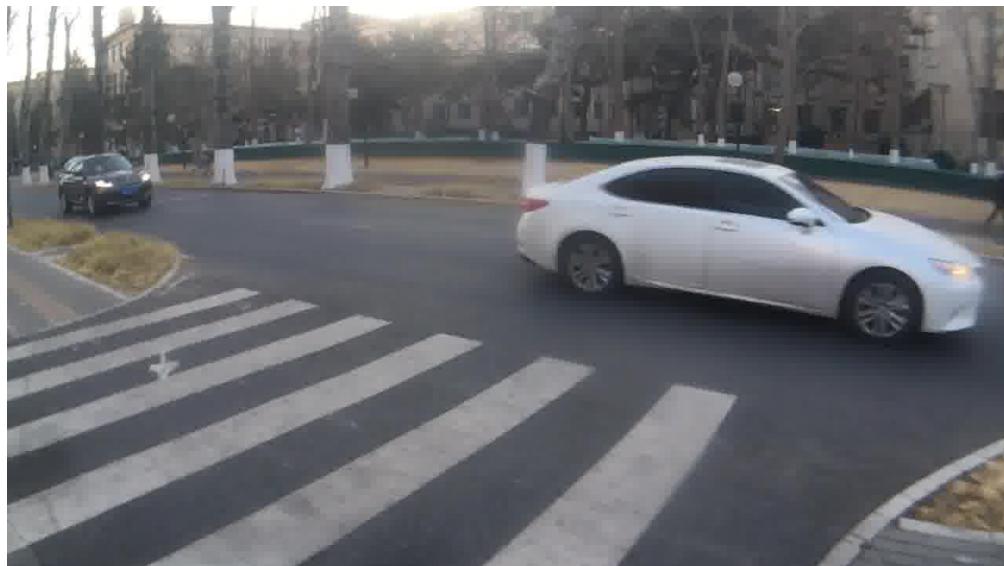
Transformed Result

Given an input video...

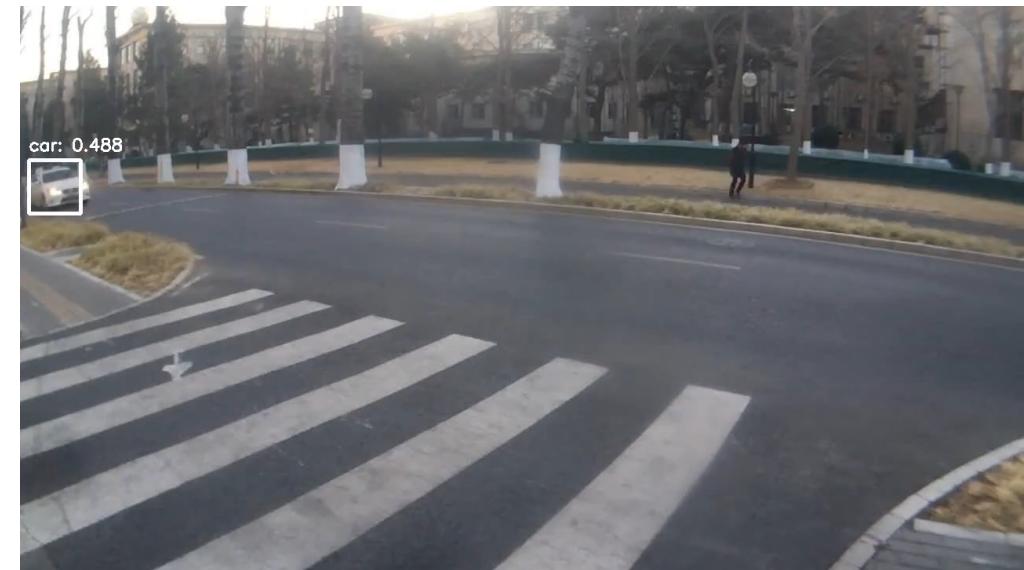


Self-supervised moving vehicle tracking with stereo sound. Gan et.al. ICCV 19

Given an input video...

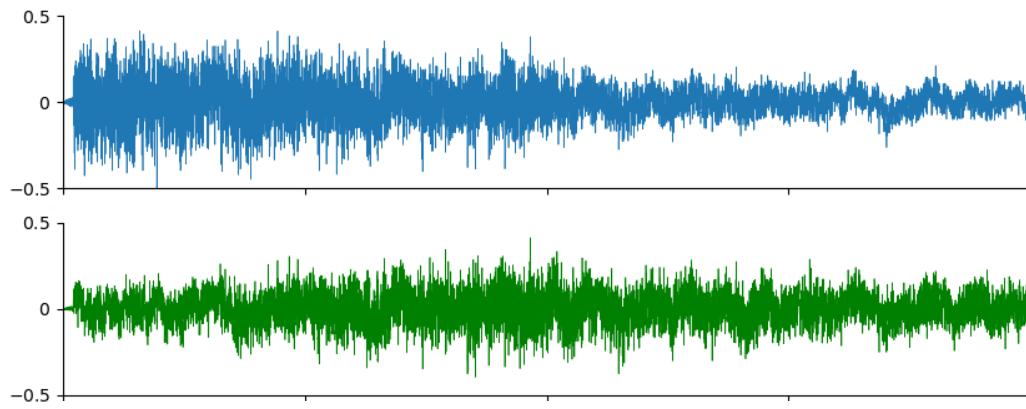


Visual detection network can track the vehicles using visual input.

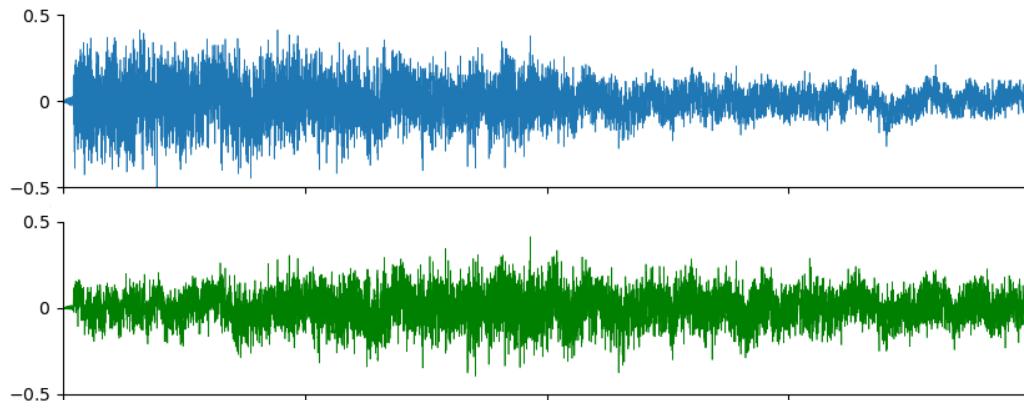


Visual tracking

What if given a piece of input stereo sound only?



What if given a piece of input stereo sound only?



Where are the vehicles?

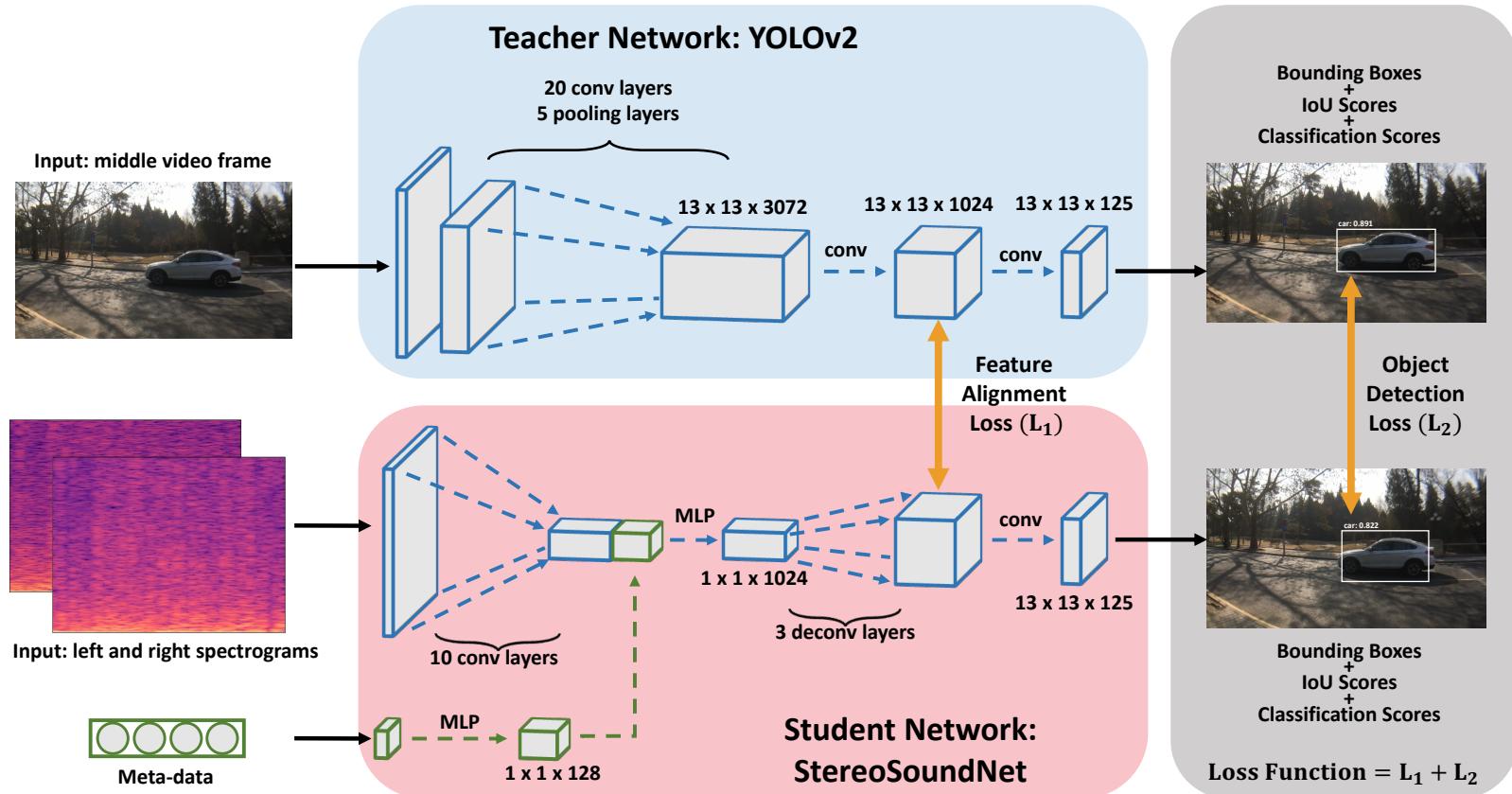
Applications

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



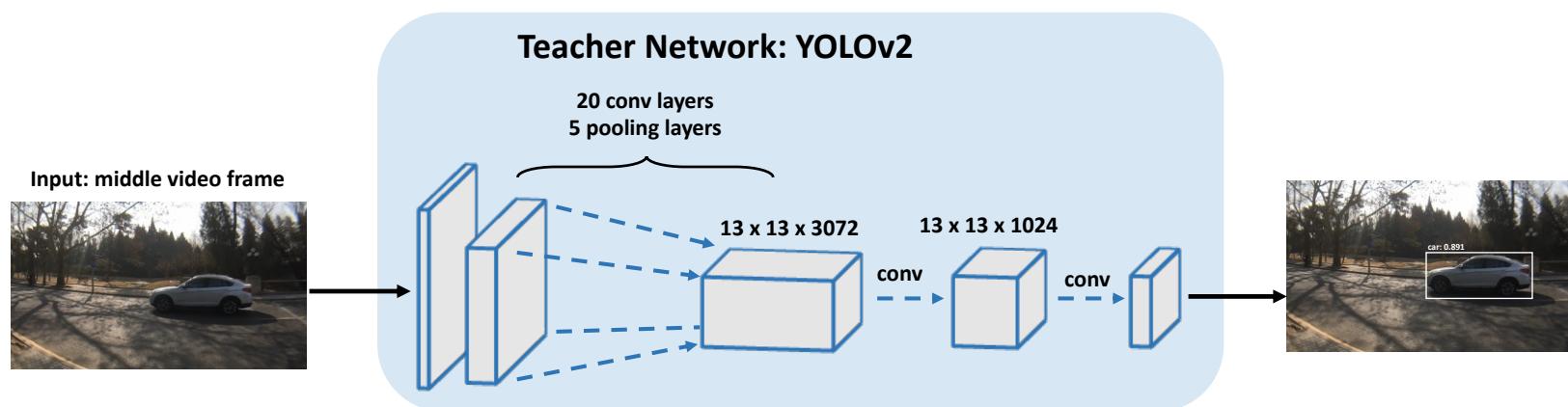
Our Methods

- Teacher-student alignment



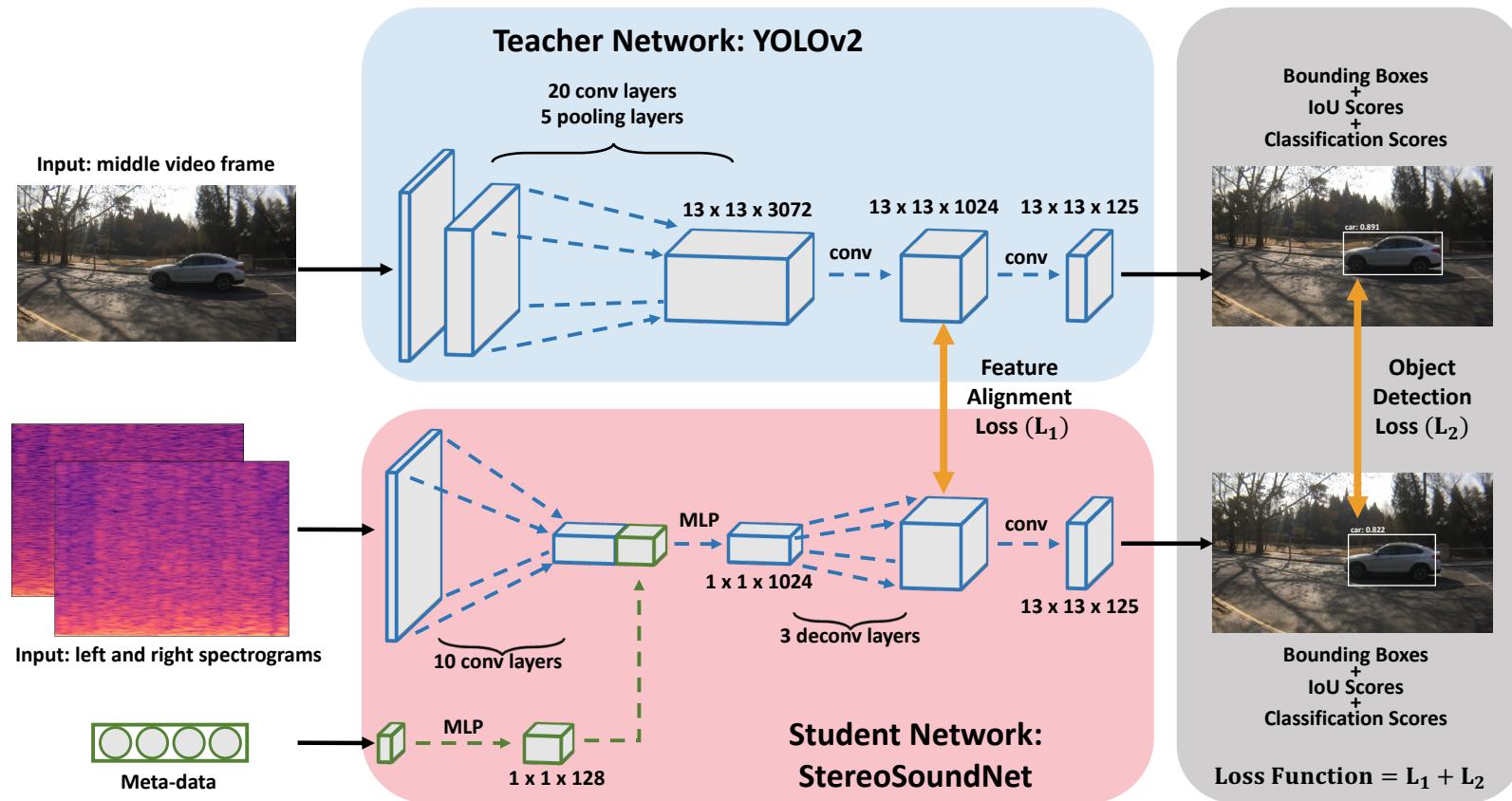
Our Methods

- 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



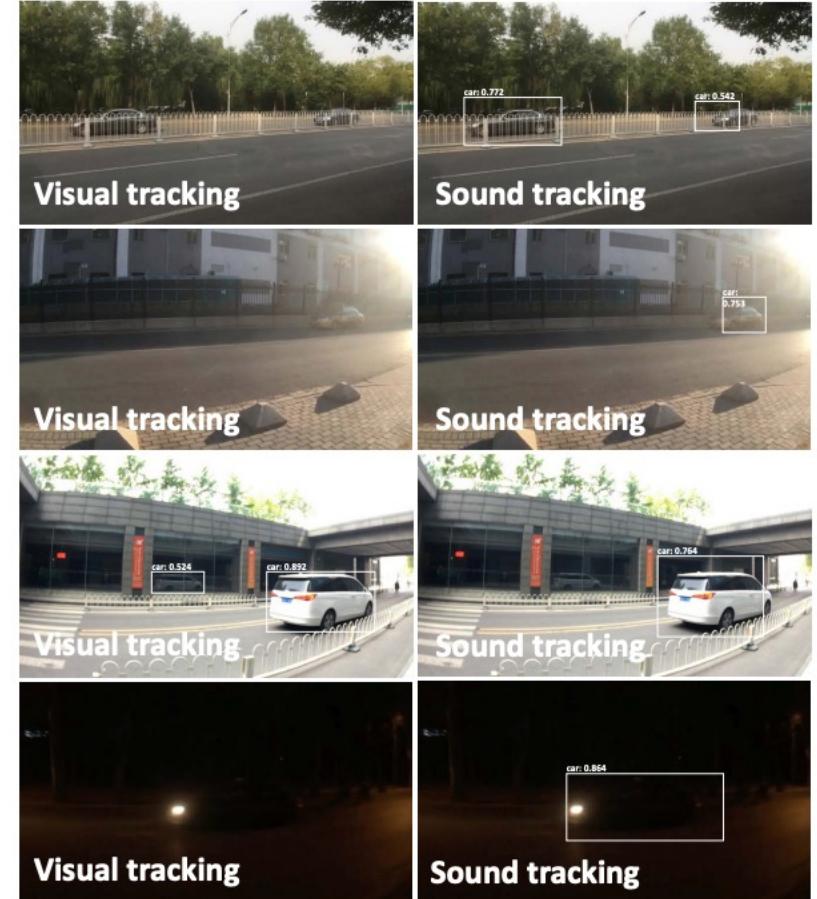
Datasets

- 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



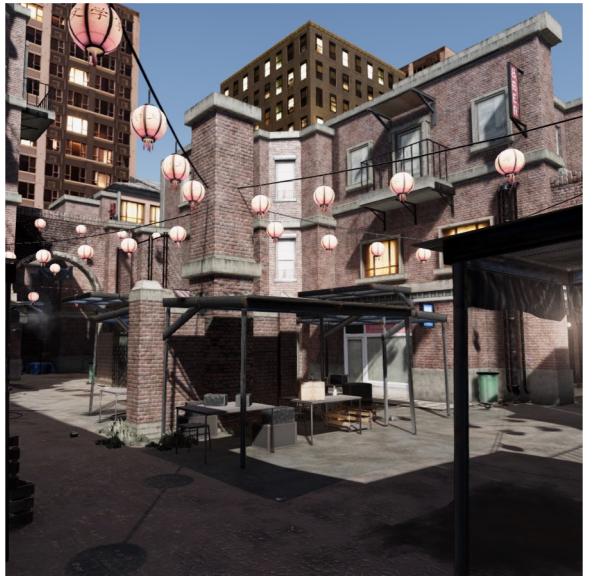
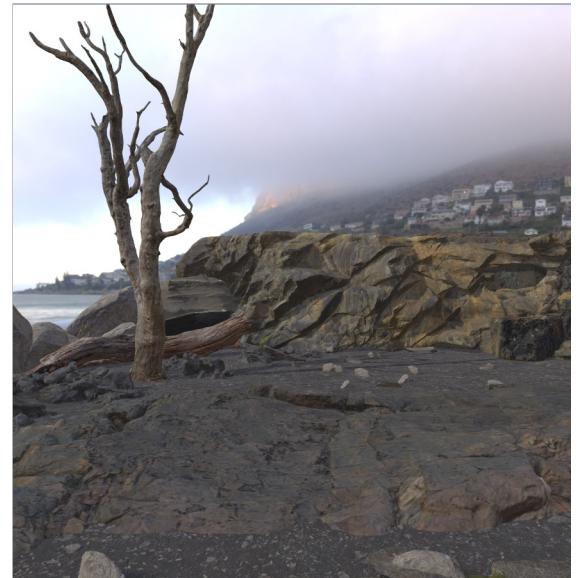
Sound Tracking

Visual Tracking

TheeDWorld: 3D Virtual World



Vision: Photo-Realistic Rendering



Audio: Physics-Triggered Sound



Physics Simulation



Rigid-body



Soft-body

Agent Interaction

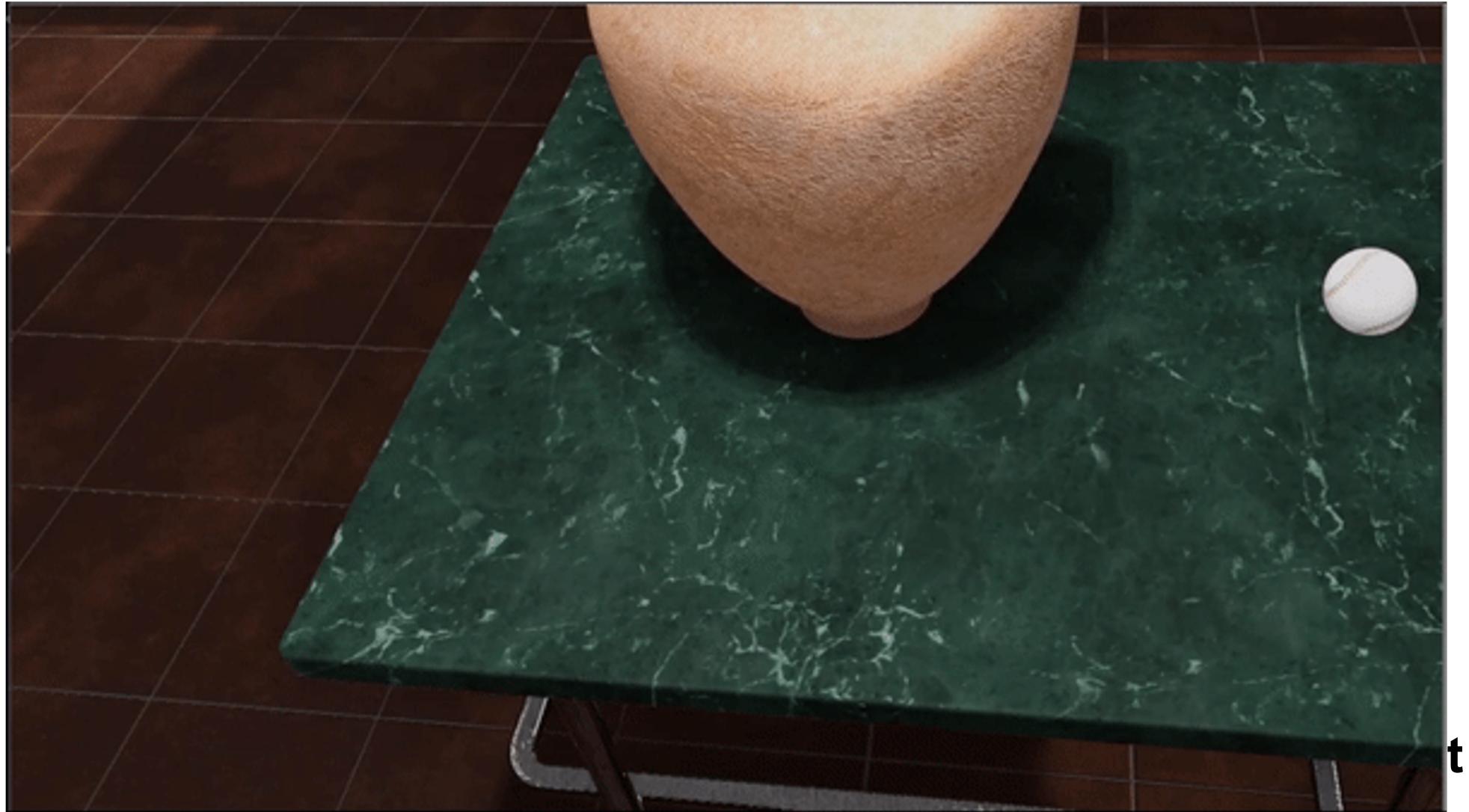


Robot



Humanoid Avatar

Object Interaction in VR

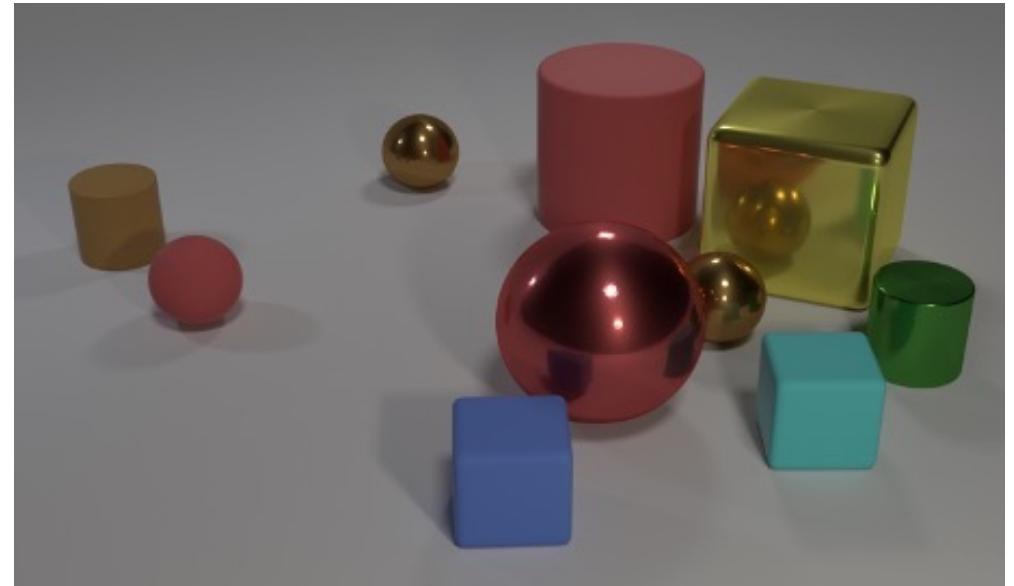


Task: Visual Reasoning



Q: What color is the fire hydrant?

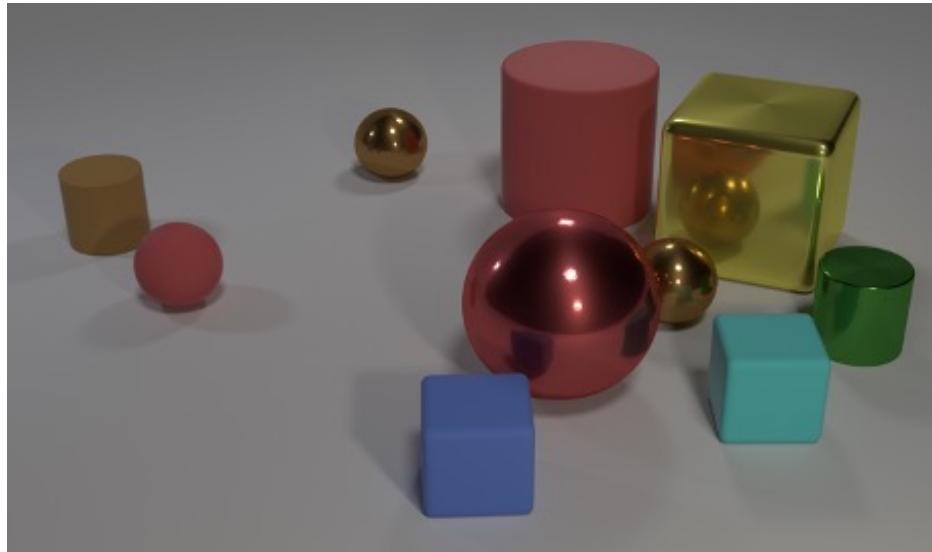
A: Yellow



Q: Are there an equal number of large things and metal spheres?

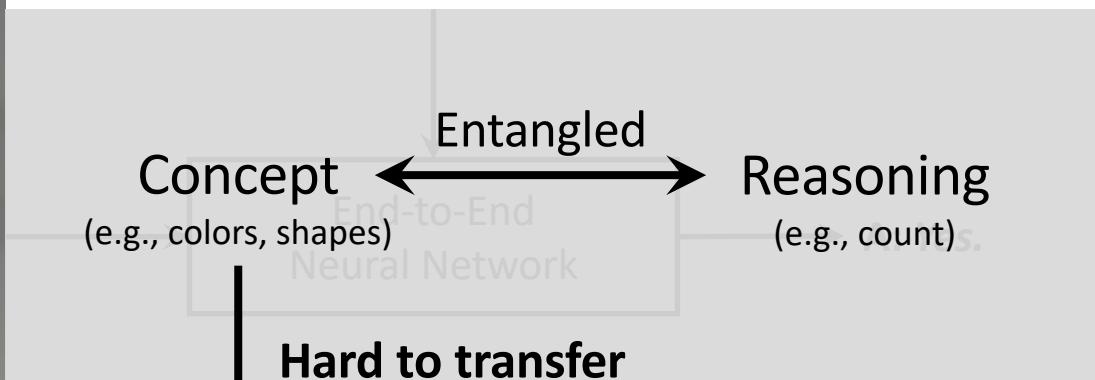
A: Yes

Prior Work: End-to-End Visual Reasoning



Visual Question Answering

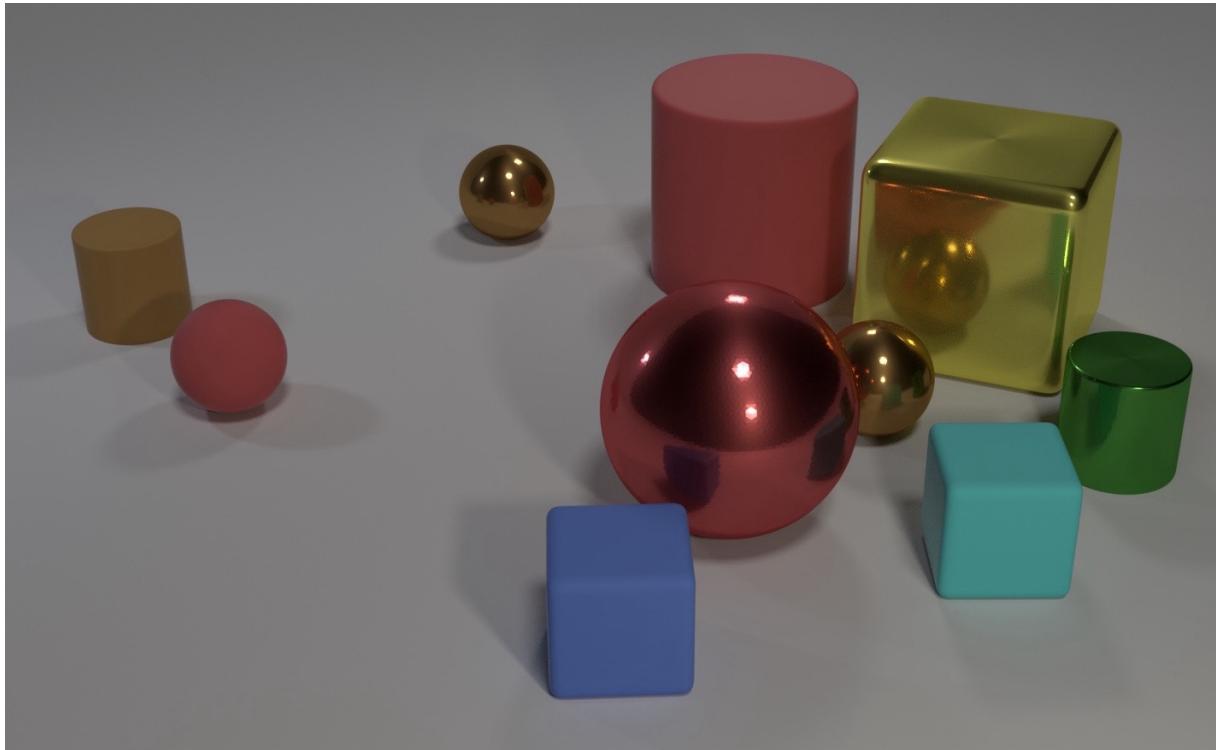
Q: Are there an equal number of **large things** and **metal spheres**?



Agrawal et al. VQA, 2015. Johnson et al. CLEVR 2017.

Johnson et al. CLEVR 2017, Andreas et al. NMN, 2016. Johnson et al. IEP, 2017. Perez et al. FiLM, 2018. Hudson & Manning. MAC, 2018. Hu et al. Stack-NMN, 2018. Mascharka et al. TbD, 2018.

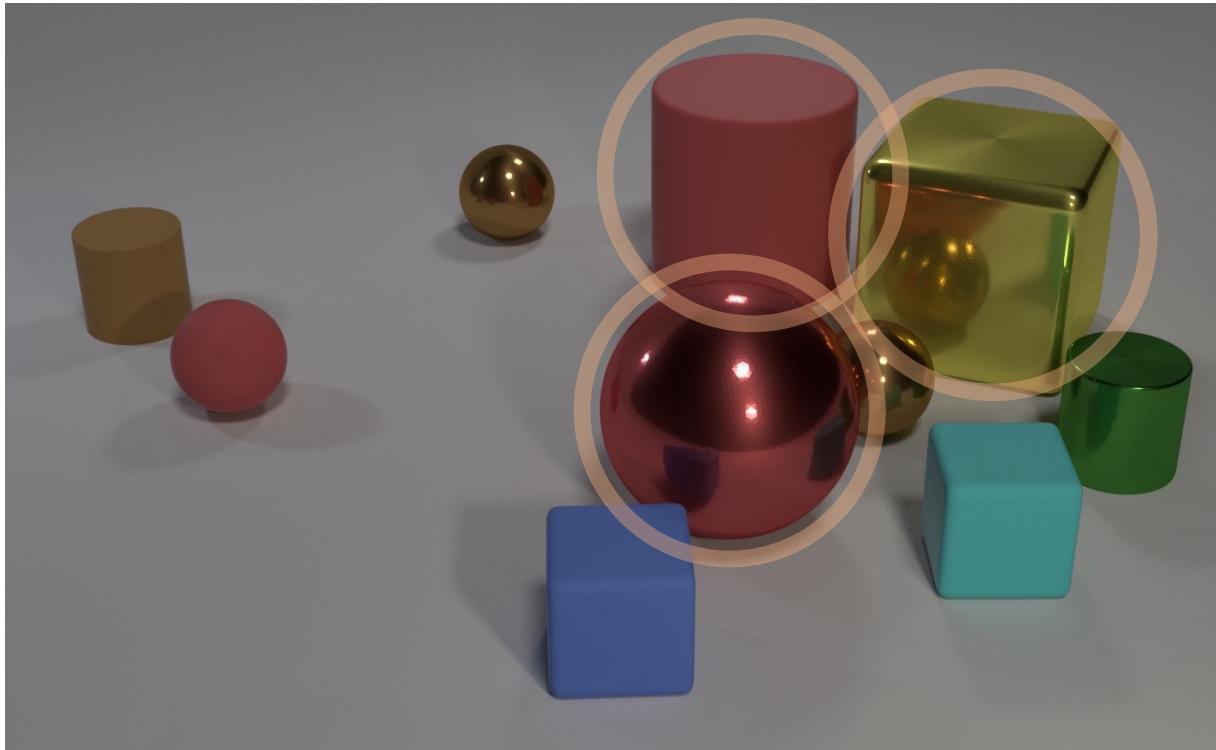
How Do Human Reason From A Visual Scene?



Question: *Are there an equal number of large things and metal spheres?*

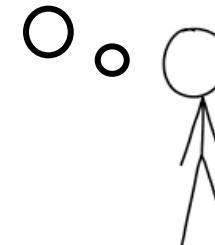


How Do Human Reason From A Visual Scene?

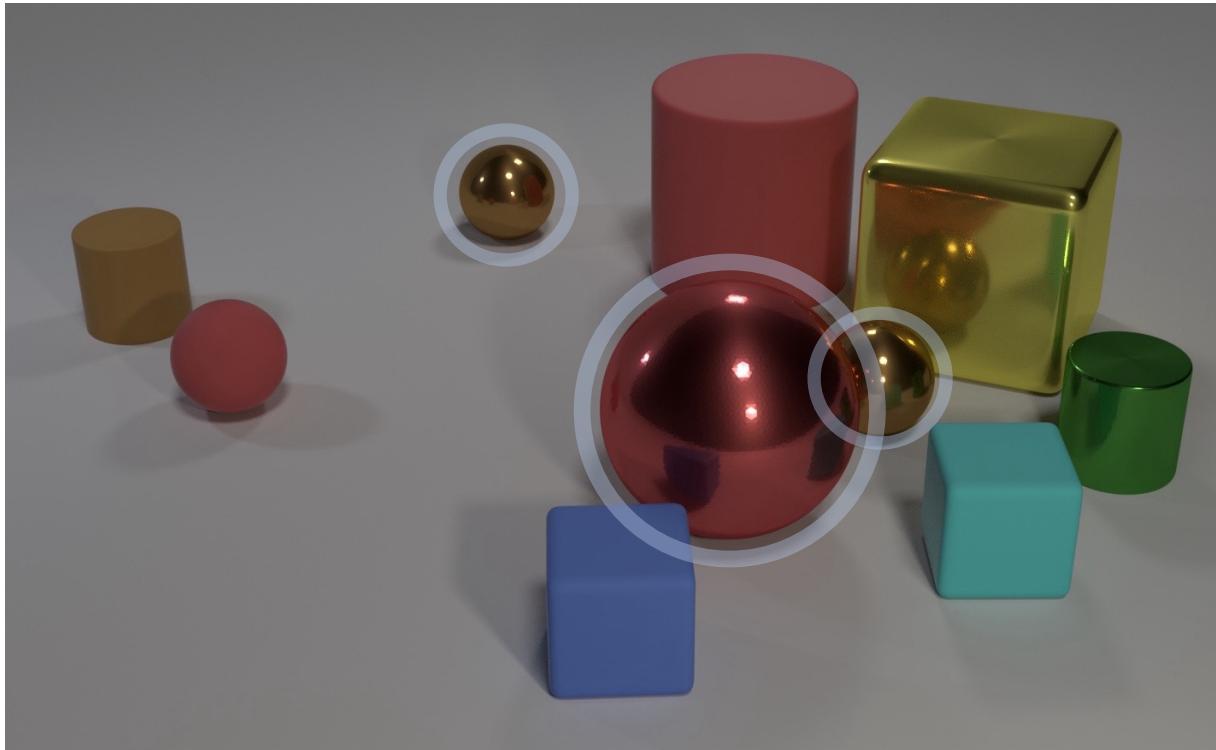


Question: Are there an equal number of *large things* and metal spheres?

3 large
things!



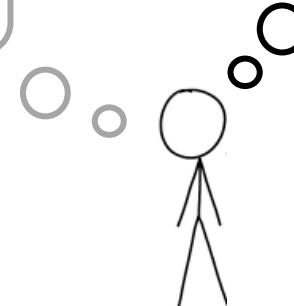
How Do Human Reason From A Visual Scene?



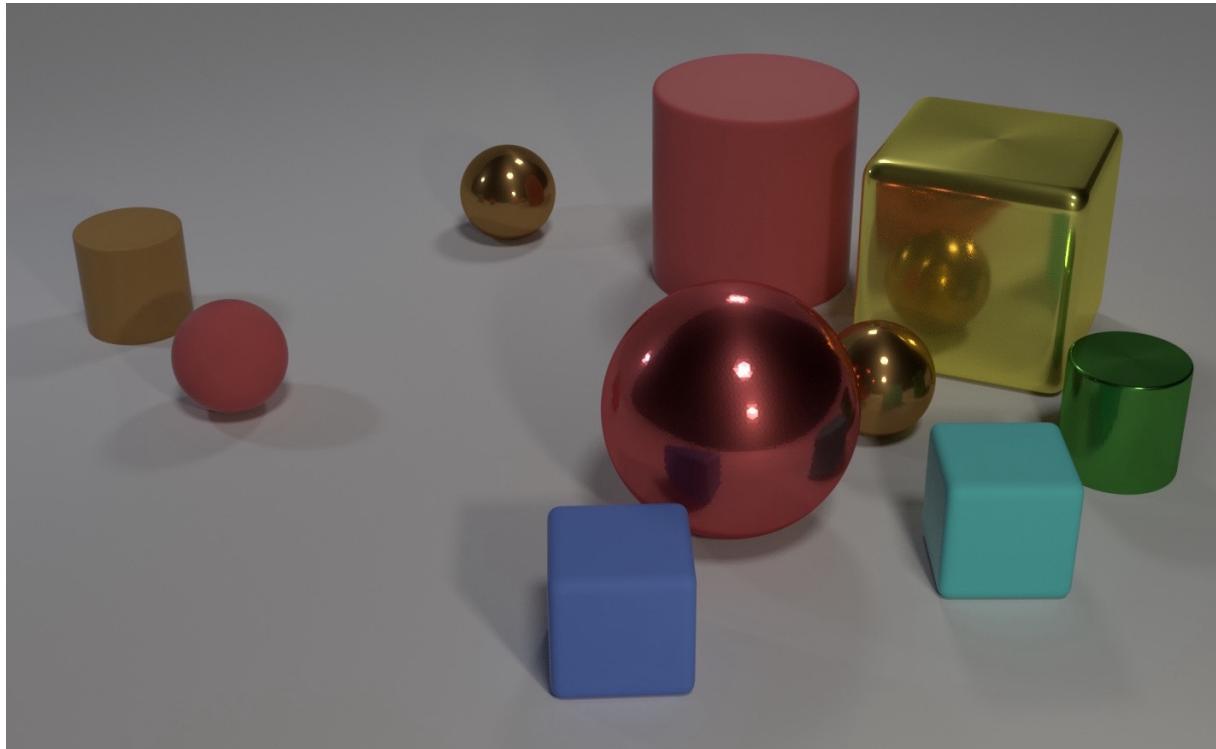
Question: Are there an equal number of large things and metal spheres?

3 large things!

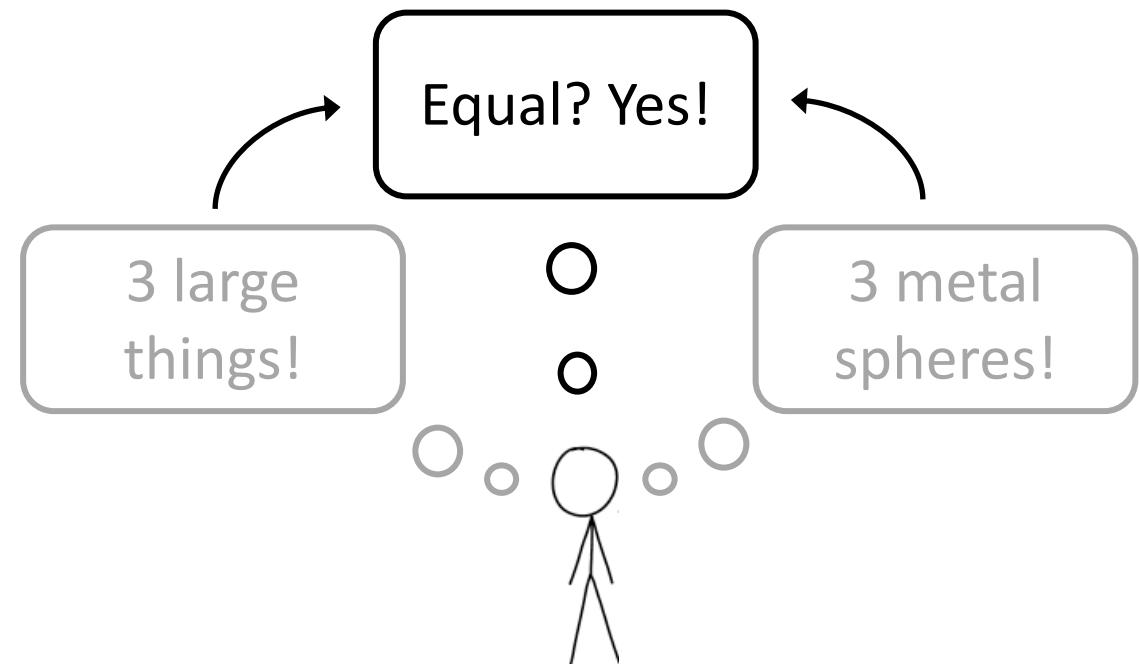
3 metal spheres!



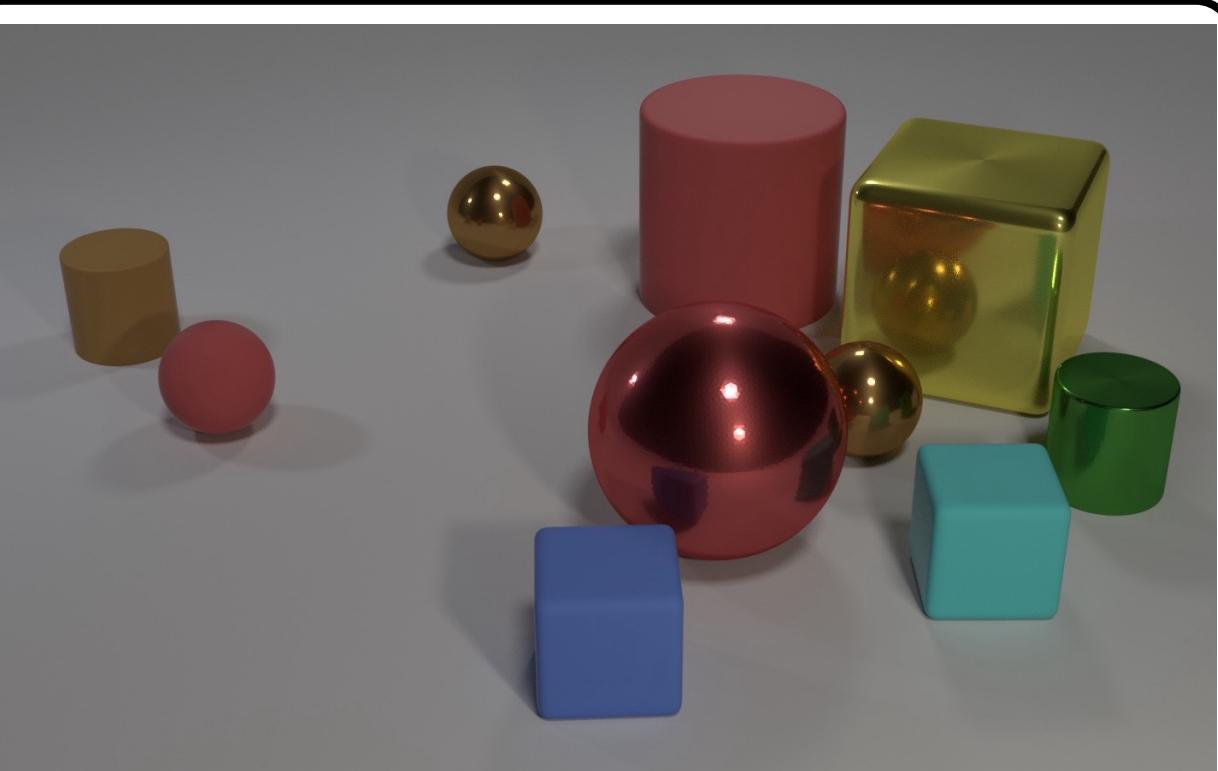
How do human reason from a visual scene?



Question: Are there an *equal number* of large things and metal spheres?



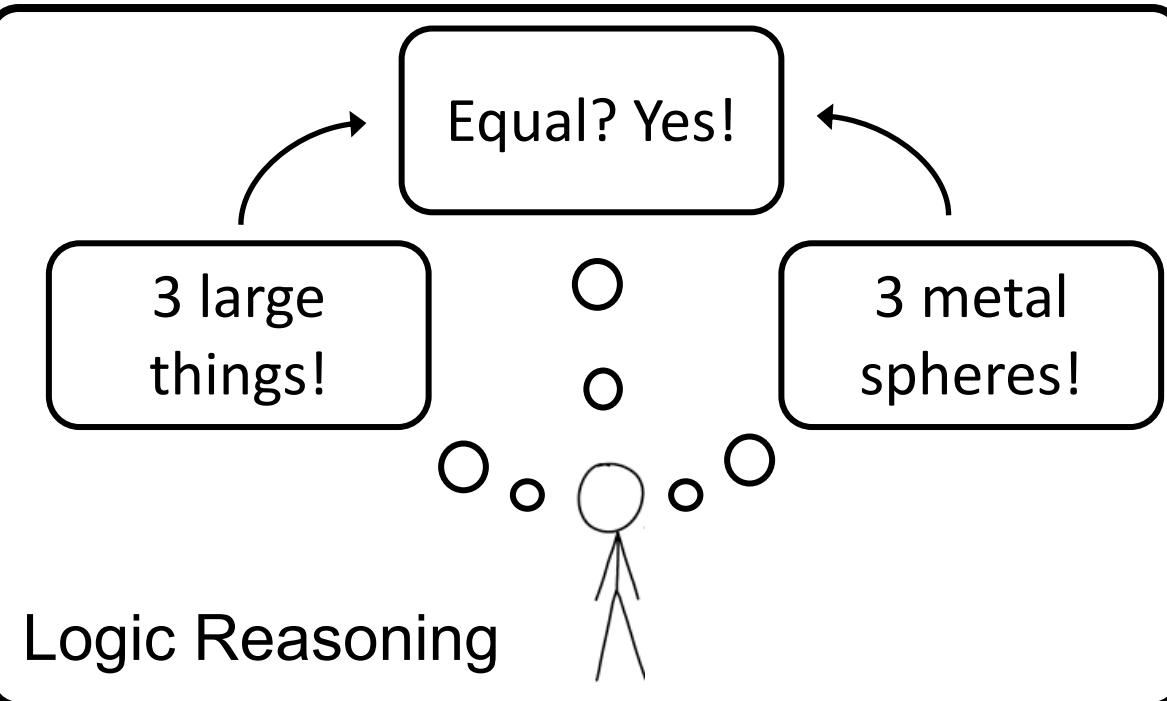
How Do Human Reason From A Visual Scene?



Visual Perception

Question Understanding

Question: *Are there an equal number of large things and metal spheres?*



Incorporate Concepts and Symbolic Programs



Neural networks parse images in **symbolic** concepts

ID	Size	Shape	Material	Color	x	y	z
1	Large	Cube	Metal	Green	-0.45	-1.10	0.35
2	Small	Cube	Metal	Green	-0.45	1.31	0.35
3	Large	Cube	Metal	Green	1.58	-1.60	0.70

I. Neural Scene Parsing

II. Neural Question Parsing



→ 1. filter_shape(scene, cylinder)

Neural networks parse questions into **symbolic** functional program.



→ 5. count(scene)

III. Symbolic Program Execution

1. filter_cylinder
2. relate_behind

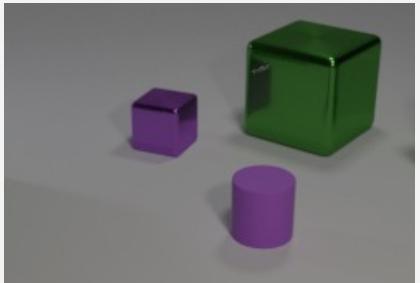
3. filter_cube
4. filter_large

5. count

Executing programs on the **symbolic** space.

3 Large Cube ...

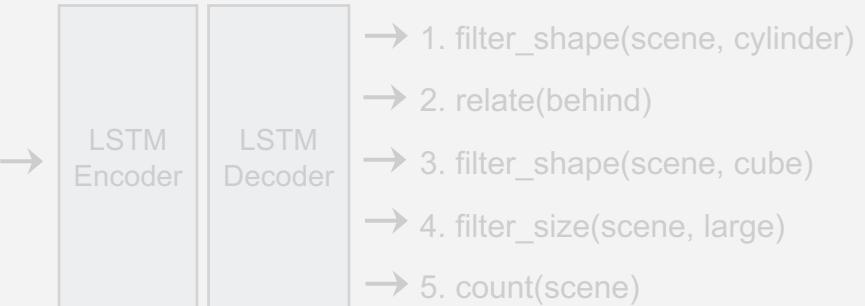
Neural Scene Parsing



I. Neural Scene Parsing

II. Neural Question Parsing

*How many cubes
that are behind
the cylinder are
large?*

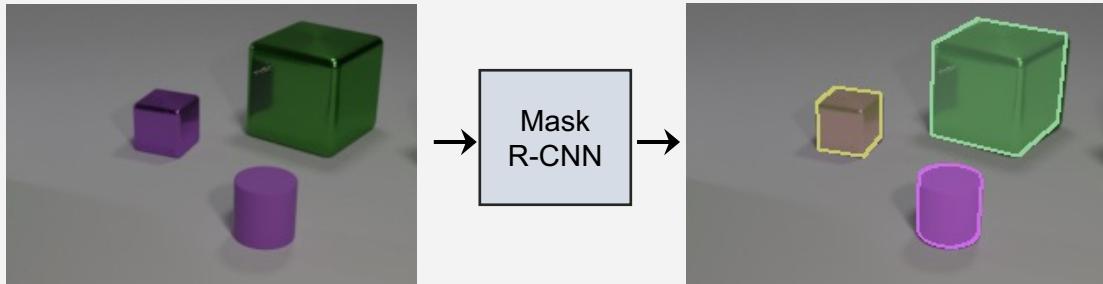


III. Symbolic Program Execution

1.	filter_cylinder	3.	filter_cube	5.	count	
2.	relate_beyond	4.	filter_large			
ID	Size	Shape	...	ID	Size	...
1	Small	Cube	...	3	Large	...
2	Small	Cylinder	...			
3	Large	Cube	...			

Answer: 1

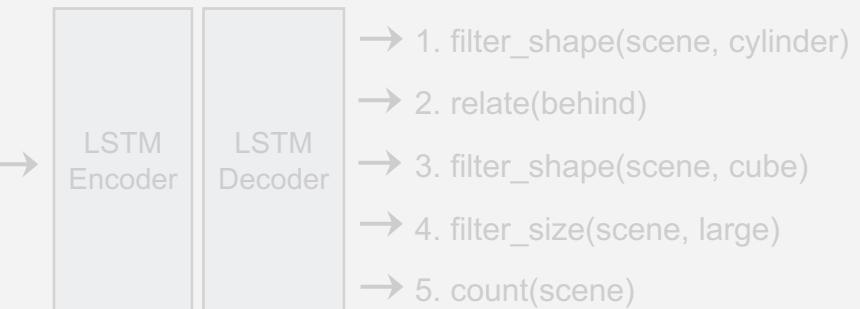
Neural Scene Parsing



I. Neural Scene Parsing

II. Neural Question Parsing

*How many cubes
that are behind
the cylinder are
large?*

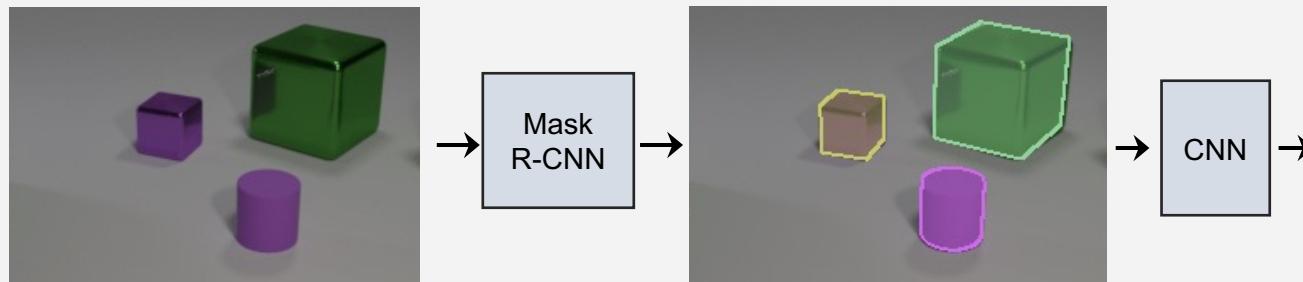


III. Symbolic Program Execution

1. filter_cylinder	3. filter_cube	5. count	
2. relate_beyond	4. filter_large		
ID	Size	Shape	...
1	Small	Cube	...
2	Small	Cylinder	...
3	Large	Cube	...

Answer: 1

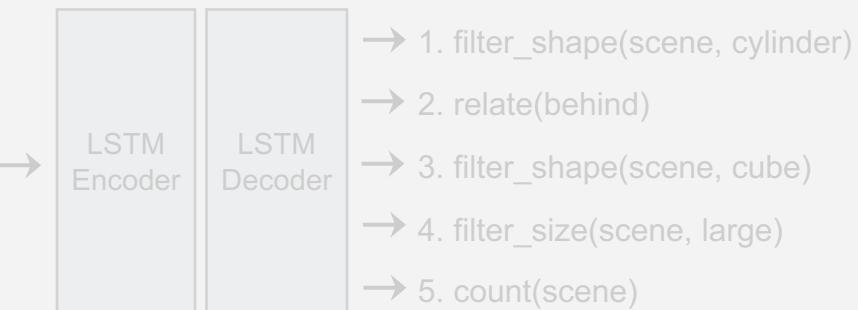
Neural Scene Parsing



I. Neural Scene Parsing

II. Neural Question Parsing

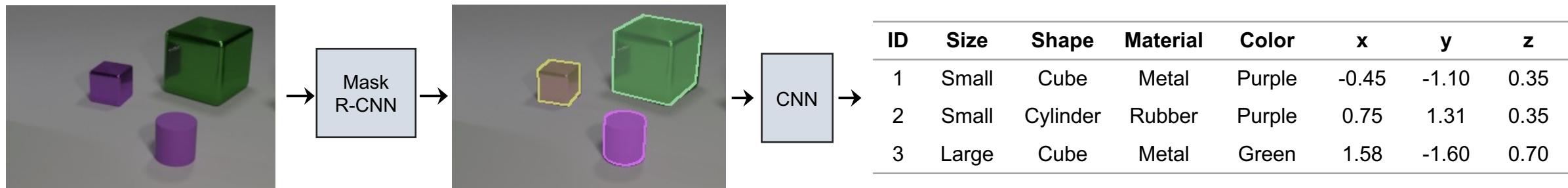
*How many cubes
that are behind
the cylinder are
large?*



III. Symbolic Program Execution



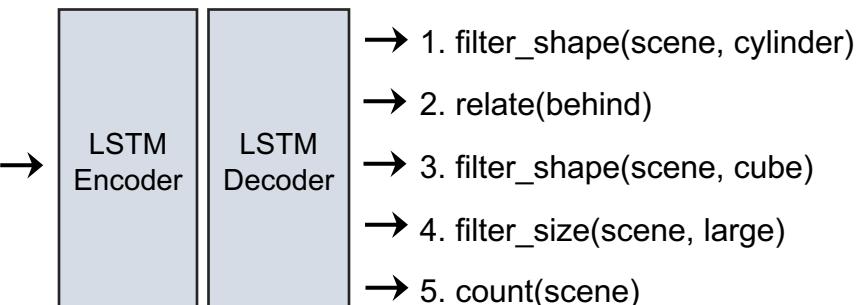
Neural Question Parsing



I. Neural Scene Parsing

II. Neural Question Parsing

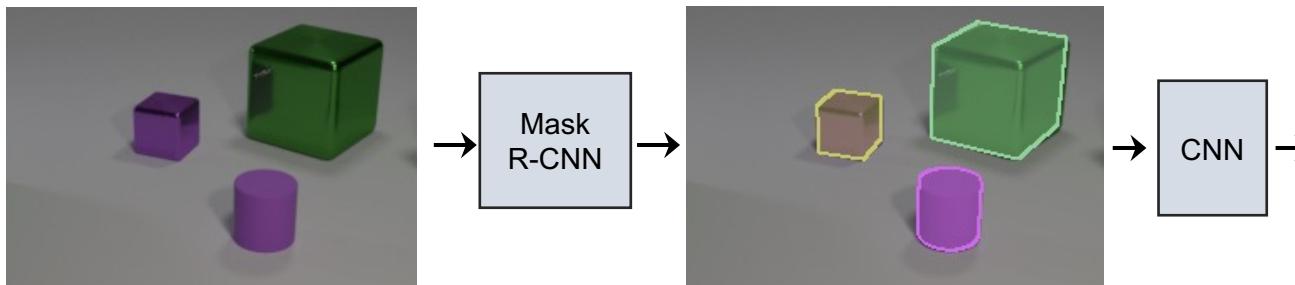
How many cubes
that are behind
the cylinder are
large?



III. Symbolic Program Execution



Symbolic Reasoning



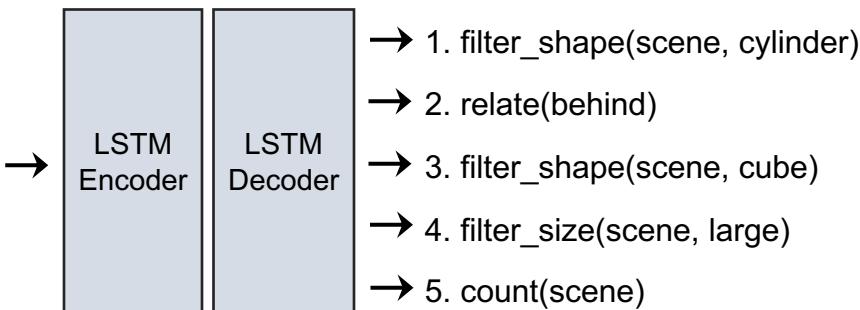
ID	Size	Shape	Material	Color	x	y	z
1	Small	Cube	Metal	Purple	-0.45	-1.10	0.35
2	Small	Cylinder	Rubber	Purple	0.75	1.31	0.35
3	Large	Cube	Metal	Green	1.58	-1.60	0.70

I. Neural Scene Parsing

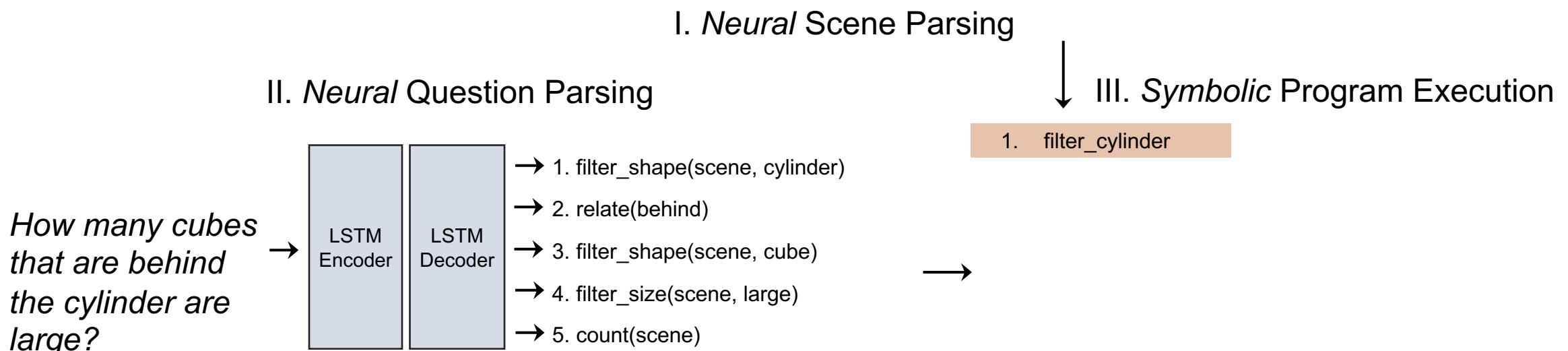
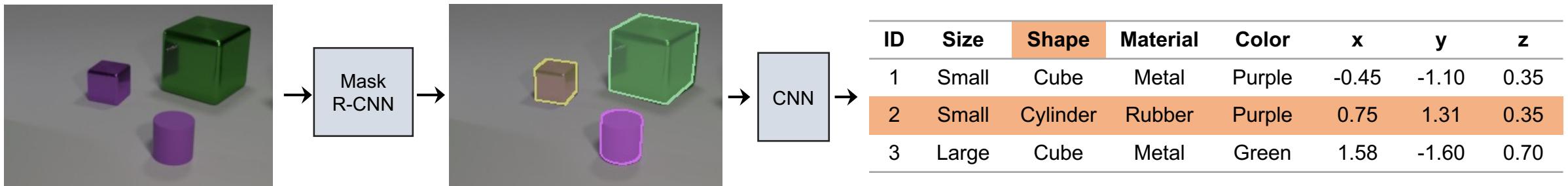
II. Neural Question Parsing

III. Symbolic Program Execution

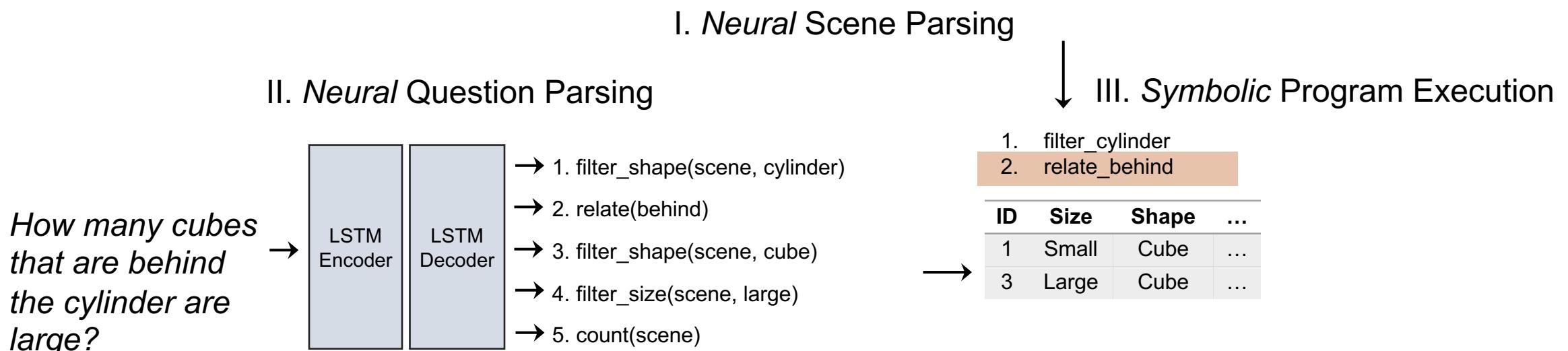
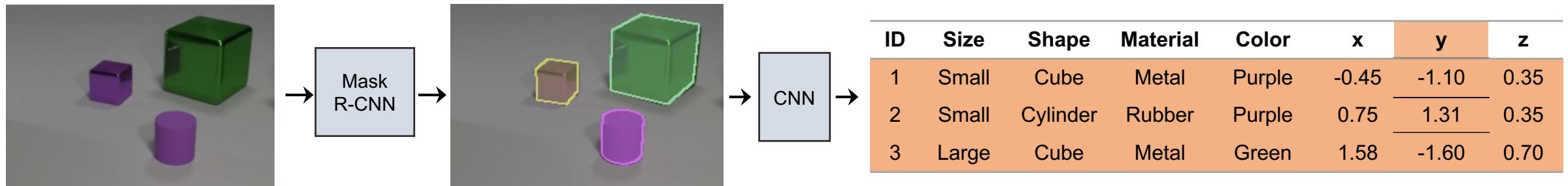
*How many cubes
that are behind
the cylinder are
large?*



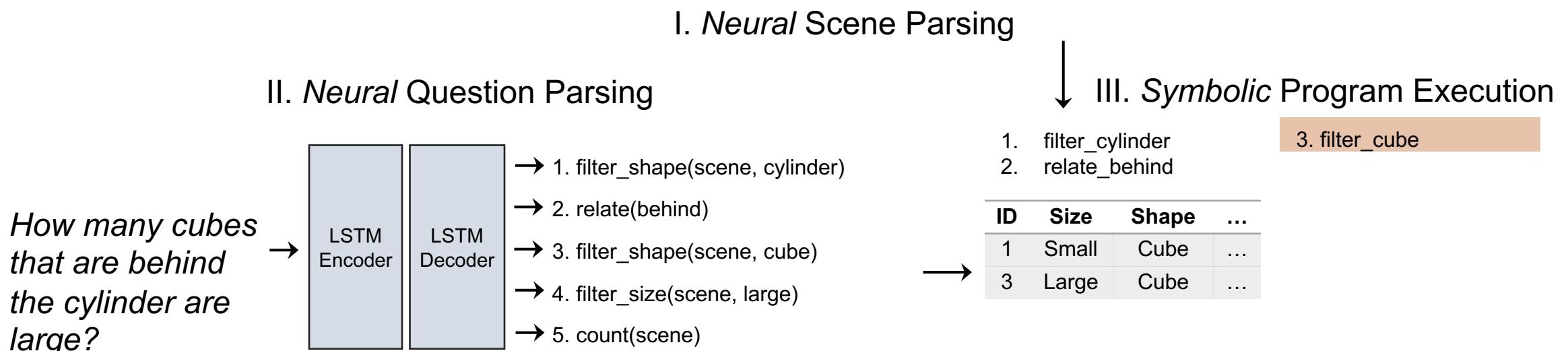
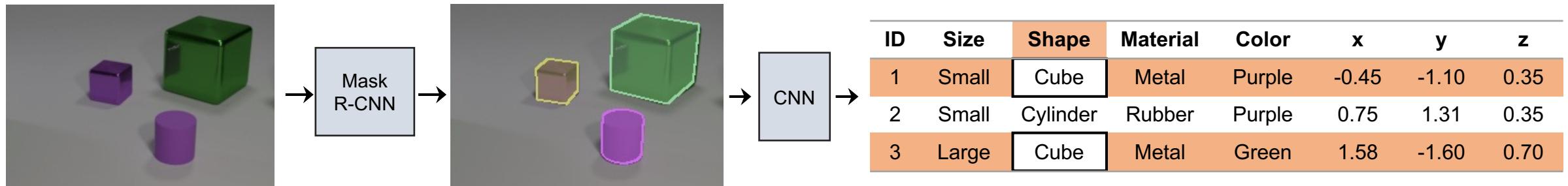
Symbolic Reasoning



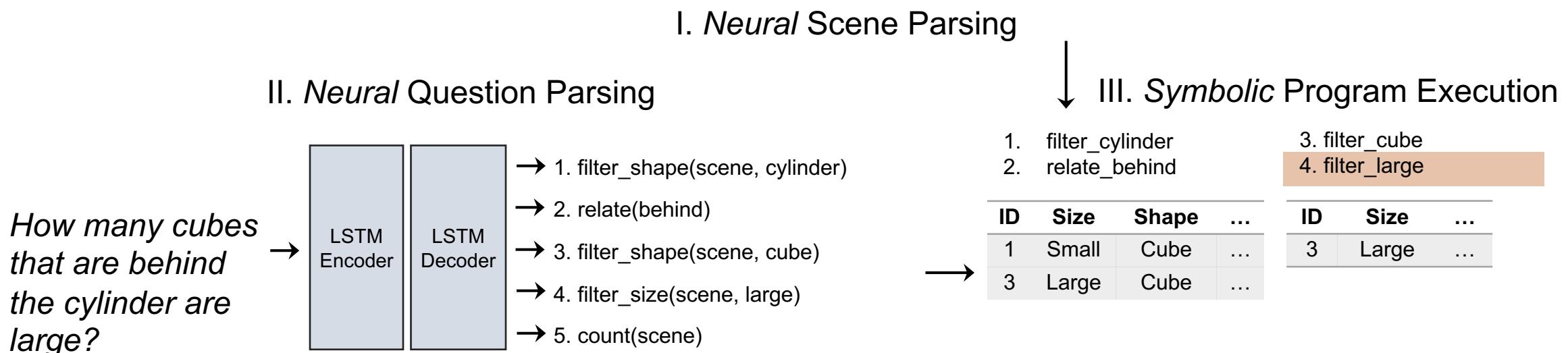
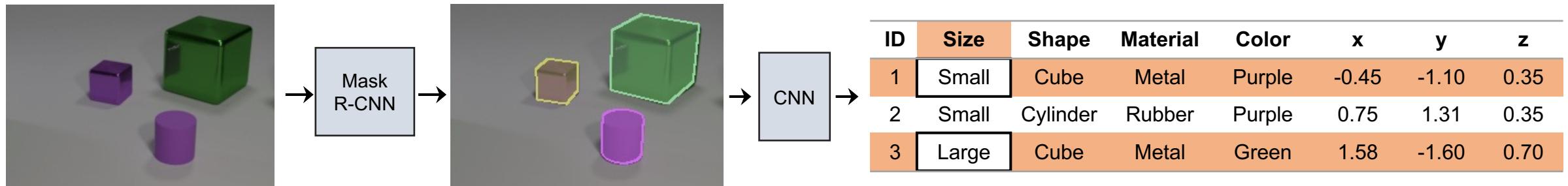
Symbolic Reasoning



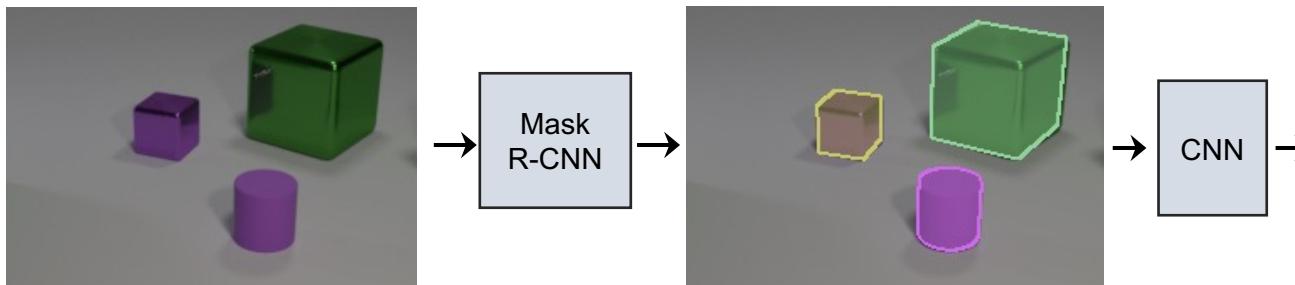
Symbolic Reasoning



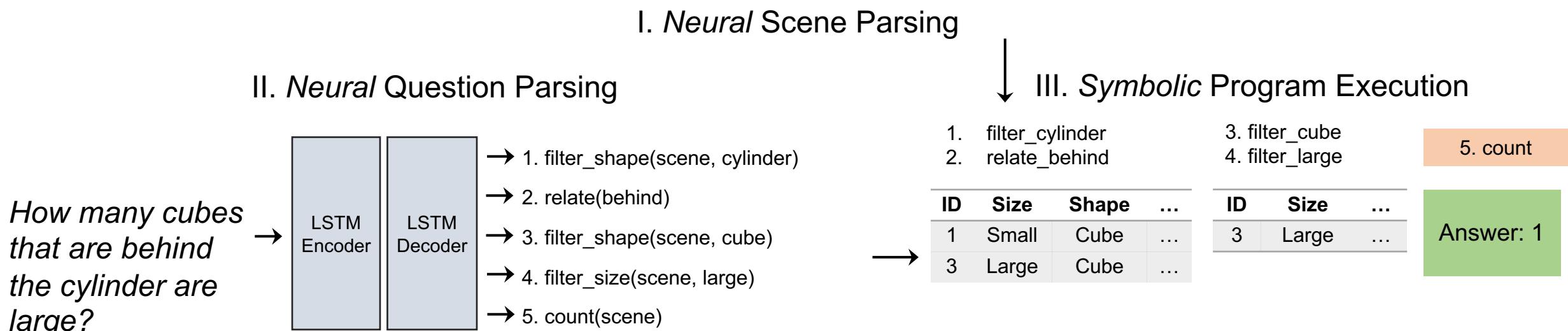
Symbolic Reasoning



Symbolic Reasoning



ID	Size	Shape	Material	Color	x	y	z
1	Small	Cube	Metal	Purple	-0.45	-1.10	0.35
2	Small	Cylinder	Rubber	Purple	0.75	1.31	0.35
3	Large	Cube	Metal	Green	1.58	-1.60	0.70



Evaluation on CLEVR

Method	Accuracy (%)
Human	92.6
RN	95.5
IEP	96.9
FiLM	97.6
MAC	98.9
NS-VQA (Ours)	99.8

