



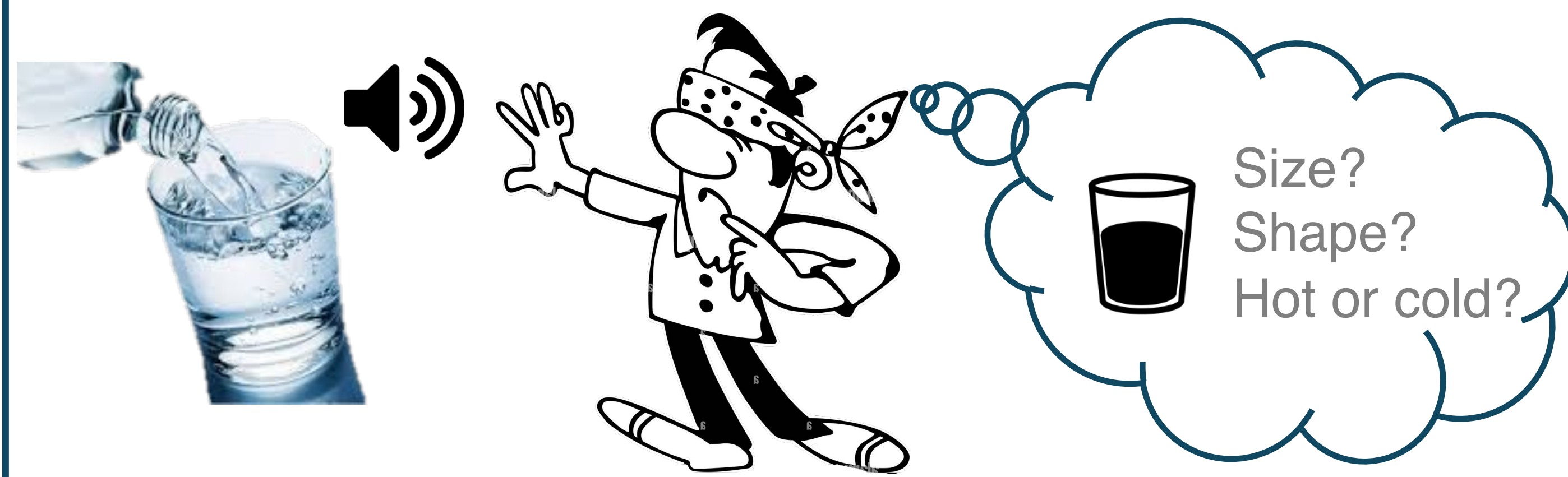
# The Sound of Water

## Inferring Physical Properties from Pouring Liquids

Piyush Bagad, Makarand Tapaswi, Cees G.M. Snoek, Andrew Zisserman

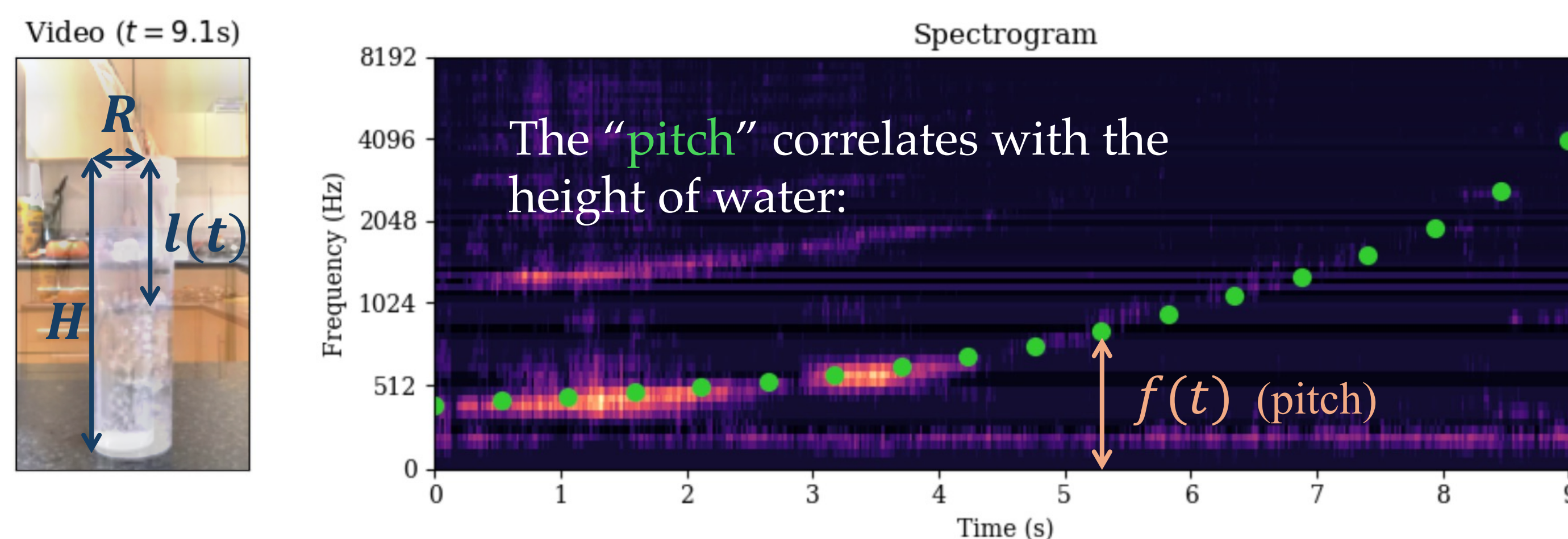


### A Remarkable Human Ability



Humans are surprisingly good at estimating physical properties merely from the sound of pouring (Cabe et al., 2000)! Can we train machines to replicate that?

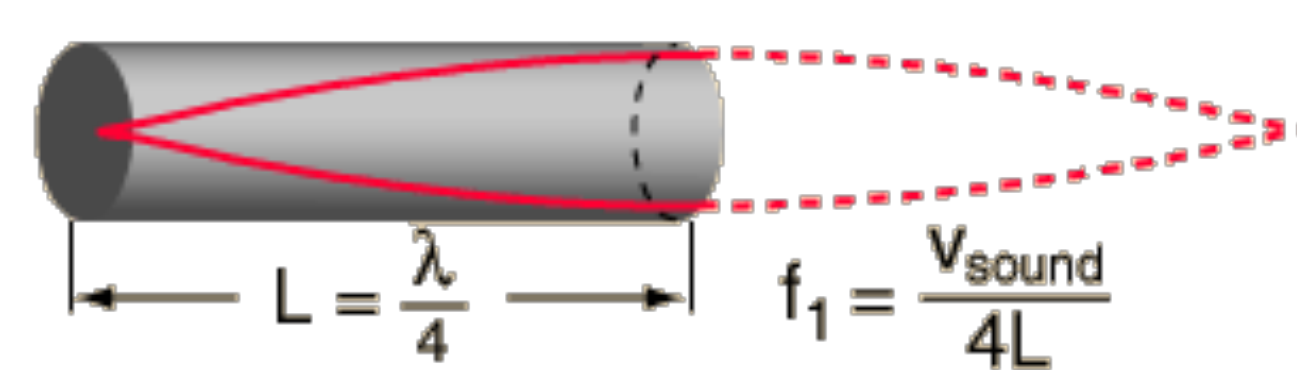
### The Physics of Pouring Sounds



Fundamental equation for the sound of pouring

$$\frac{c}{4f(t)} = l(t) + \beta R; \quad l(t) = \begin{cases} H, & t = 0 \\ 0, & t = T \end{cases}$$

- $c$  is the speed of sound in air; and  $\beta$  experimental constant
- Underlying principle is the same as that in a resonant pipe

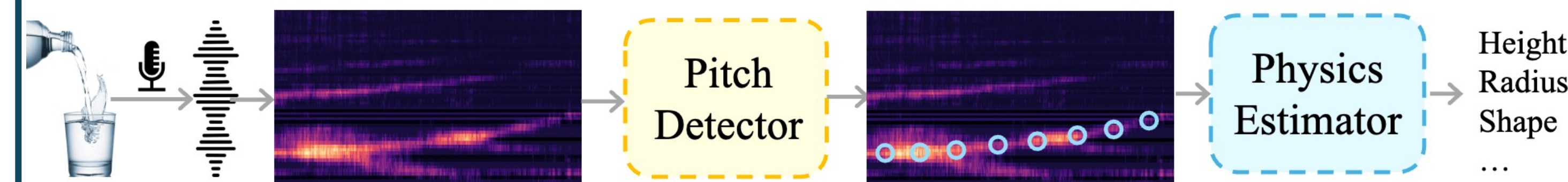


Recovering physical properties from pitch

$$l(t) = \frac{1}{4} \left( \frac{c}{f(t)} - \frac{c}{f(T)} \right); \quad H = l(0); \quad R = \frac{c}{4\beta f(t)}$$

Height  $H$  depends on accurate pitch at the start of pouring and radius  $R$  on the end of pouring

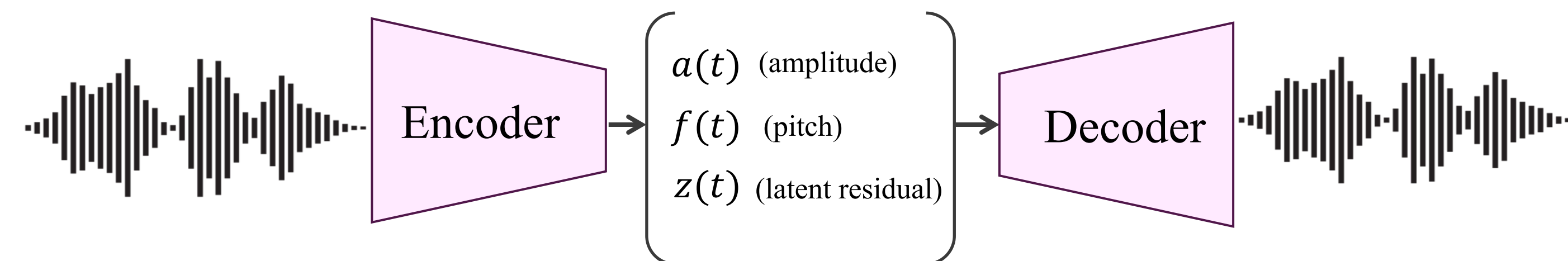
### Training Pitch Detector by Visual Co-supervision



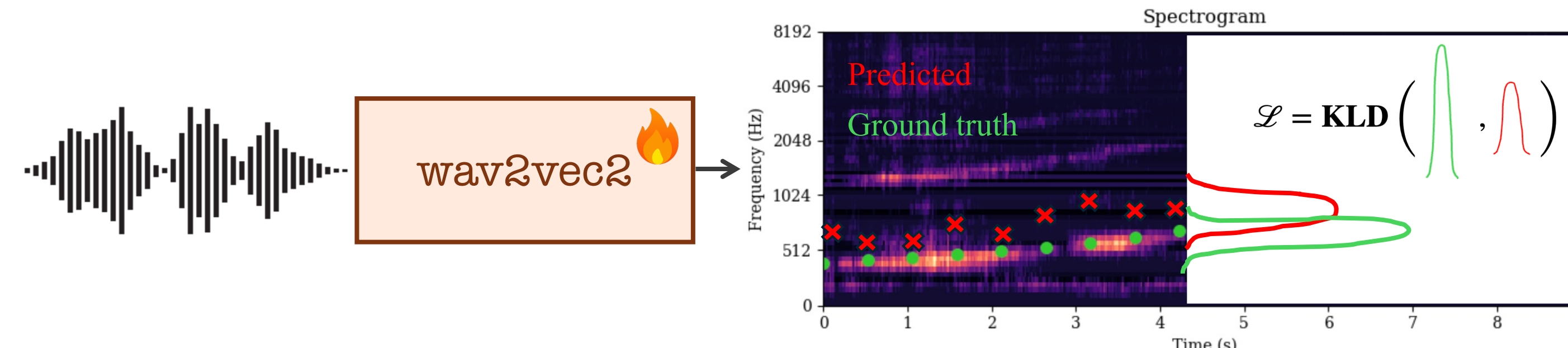
How to detect pitch in pouring sounds?

1. Simulate sounds of pouring with desired pitch profile
2. Pre-train a pitch detector network (wav2vec2) on simulated data
3. Fine-tune on real data with co-supervision from the video stream

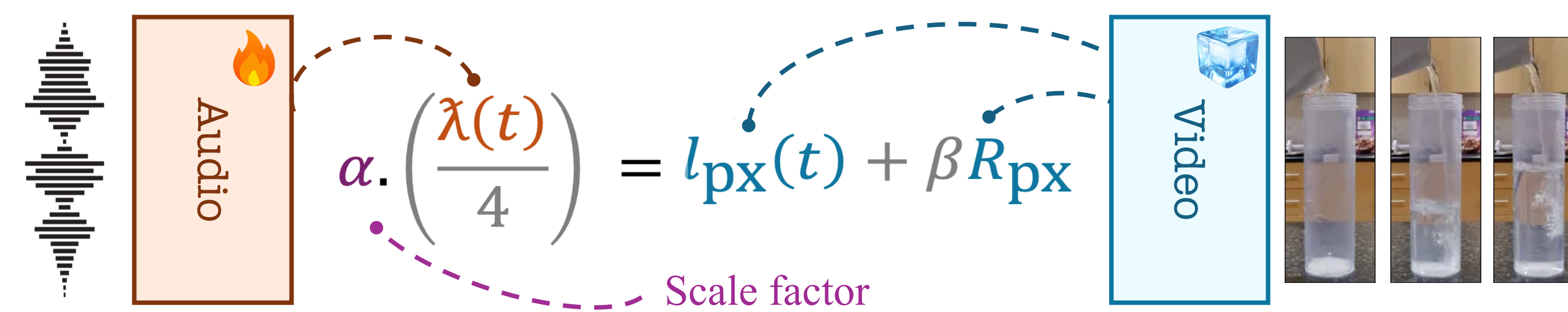
#### I. Simulate sounds of pouring



#### II. Pre-train on simulated data



#### III. Fine-tune on real data with video teacher



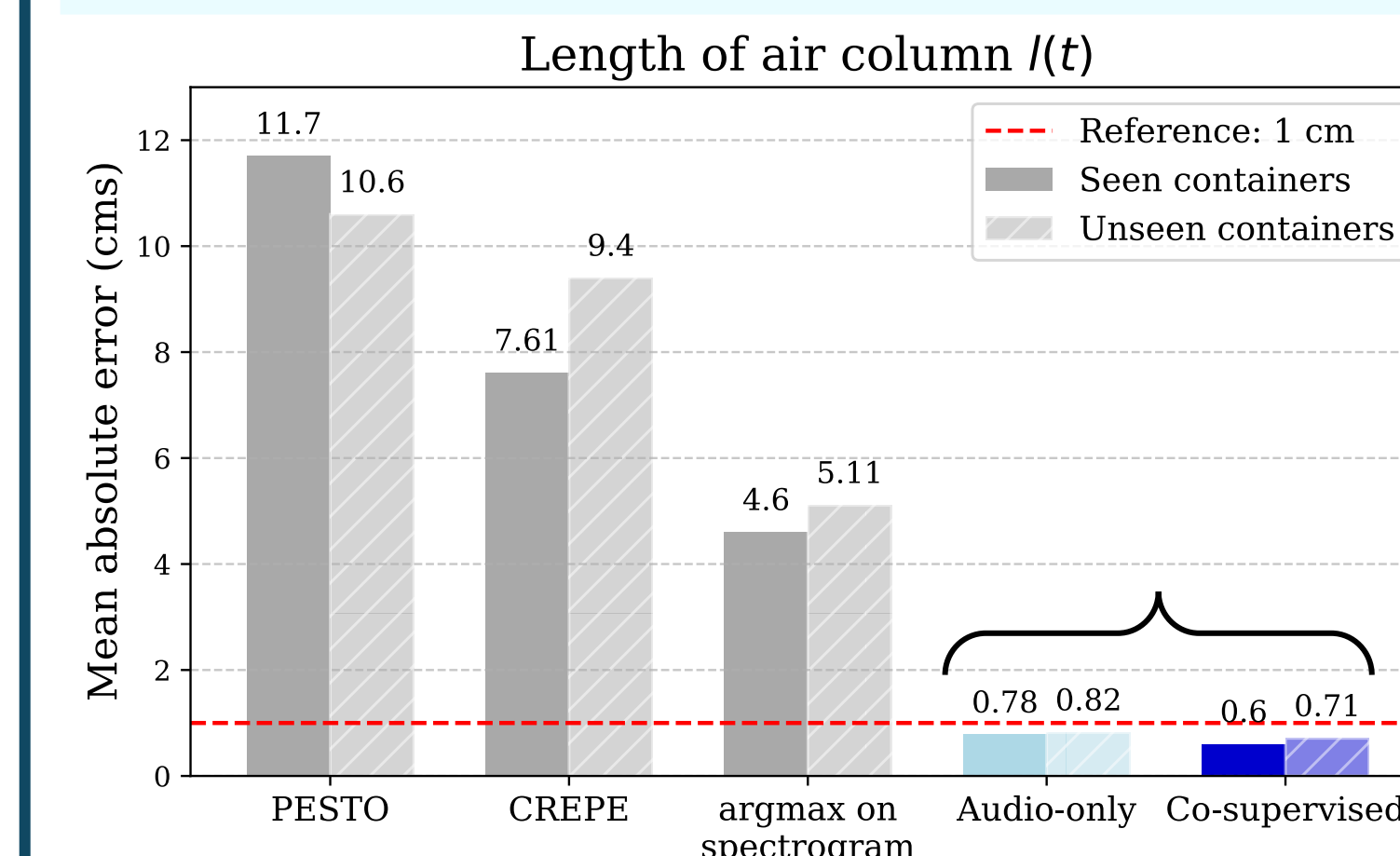
### Train and Evaluation Dataset: Sound of Water 50

For train/evaluation, 805 videos spanning 50+ containers (4 shapes, 5 materials, 2 liquids)



### Experimental Results

Achieves an error rate of < 1 cm; and co-supervision helps!



Property	Synthetic ↓	Co-supervised ↓	Δ
Length $l(t)$ (cm)	0.78	0.60	+0.18
Height $H$ (cm)	2.23	2.27	-0.04
Radius $R$ (cm)	1.62	1.39	+0.23
Flow rate $Q$ (ml/s)	25.2	22.5	+2.70
Time to fill $\tau$ (s)	1.62	1.49	+0.13

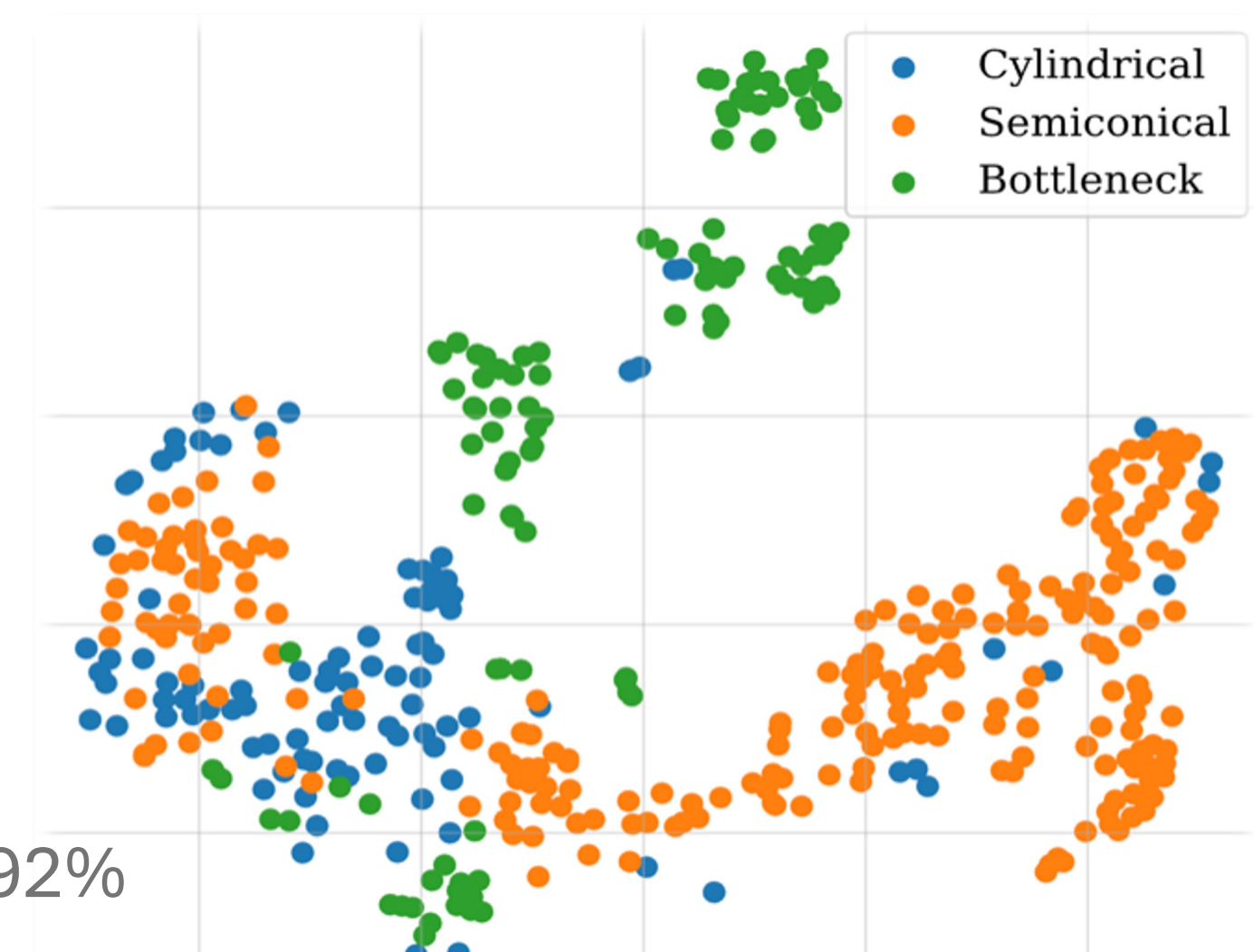
1. Our models achieve < 1cm error in estimating length of air column
2. In general, co-supervision tends to help beyond synthetic pre-training

The learned features encode liquid mass and container shape!

Samples from Wilson et al (2019)



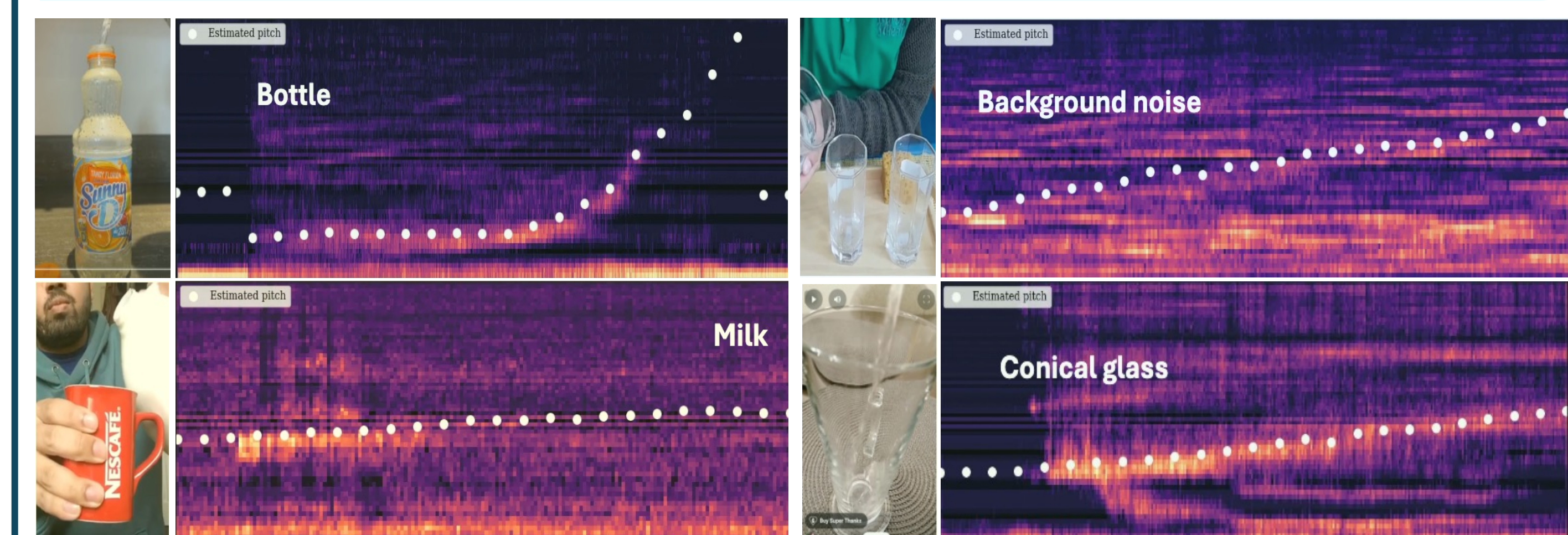
t-SNE of learned representations



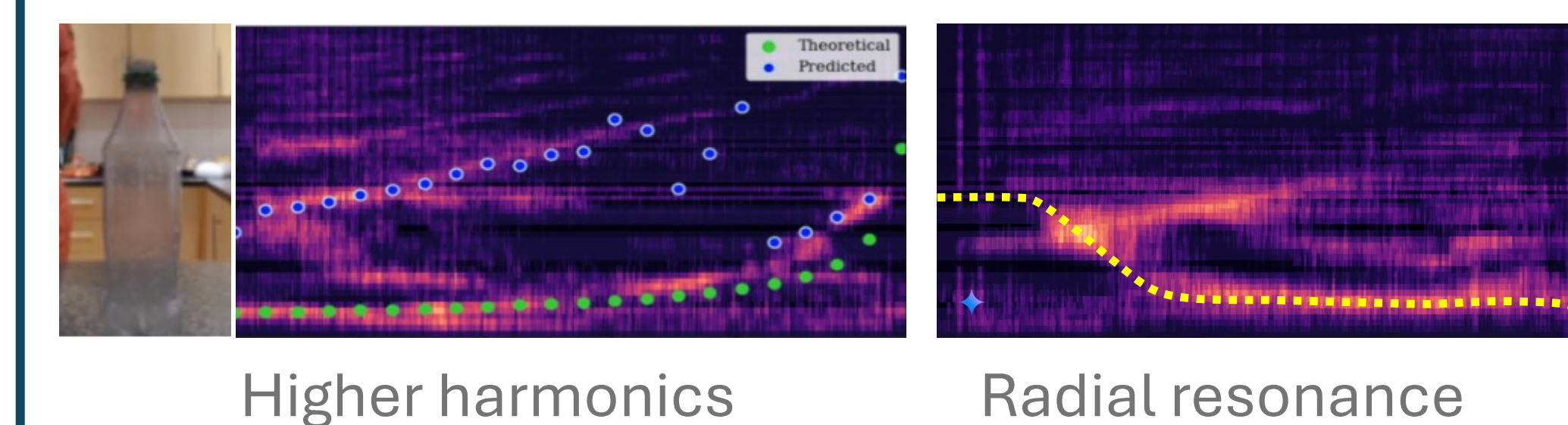
Liquid mass estimation on Wilson et al. (ICRA 2019): MAE: 34ml

Shape classifier accuracy: 92%

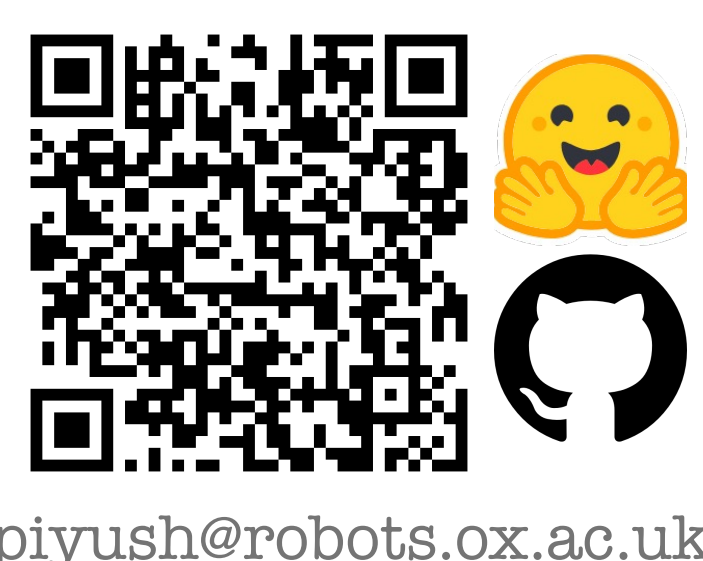
Generalization to novel container shapes, materials, liquids and even in-the-wild YouTube samples.



Failure cases and future work



Code & Models



piyush@robots.ox.ac.uk