# Self-Labelling via simultaneous clustering and representation learning

YUKI M. ASANO* MANDELA PATRICK* CHRISTIAN RUPPRECHT ANDREA VEDALDI

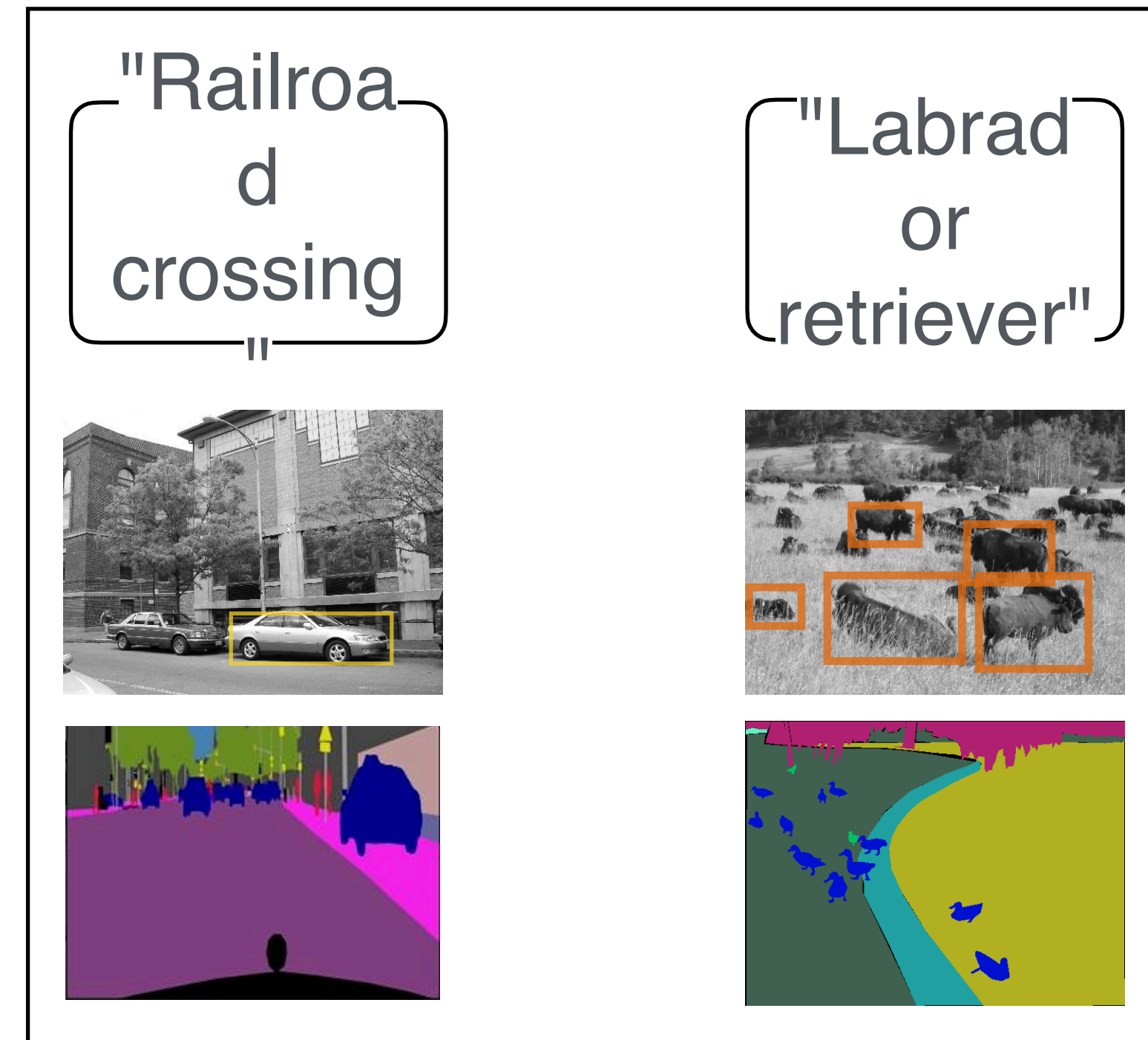yuki@robots.ox.ac.uk, @y_m_asano

UNIVERSITY OF OXFORD

VGG
UNIVERSITY OF OXFORD

FACEBOOK

# Manual annotations for the data are limiting.

Data is often cheap



But manual annotations are expensive



"Railroad crossing"

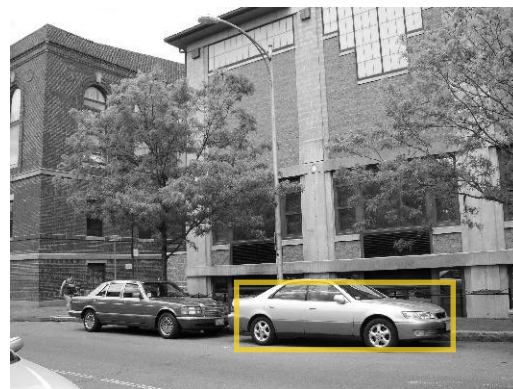"Labrador retriever"

# Replacing manual annotations by self-supervised learning.
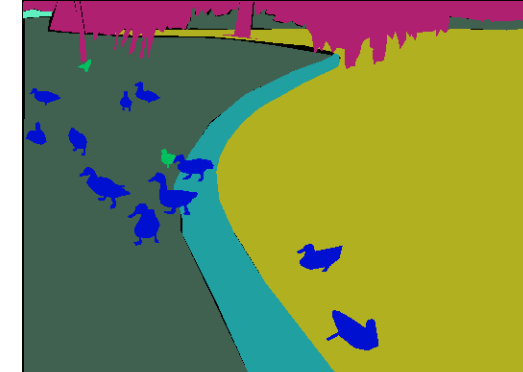


Clustering

Detection

Segmentation

**Images**

CliqueCNN (Bautista, NeurIPS'16)
DeepCluster (Caron, ECCV'18)
IIC (Ji, ICCV'19)
SeLa (Asano, ICLR'19)
SCAN (Gansbeke, ECCV'20)
and more

SSOD (Afouras arxiv'21)

MaskContrast (Gansbeke arxiv'21)

**Video**

[Sight from Sound (Owens ECCV'16)]
XDC (Alwassel NeurIPS'20)
**SeLaVi (Asano NeurIPS'20)**

Boxes:
SSOD (Afouras arxiv'21)

[Heatmaps]:
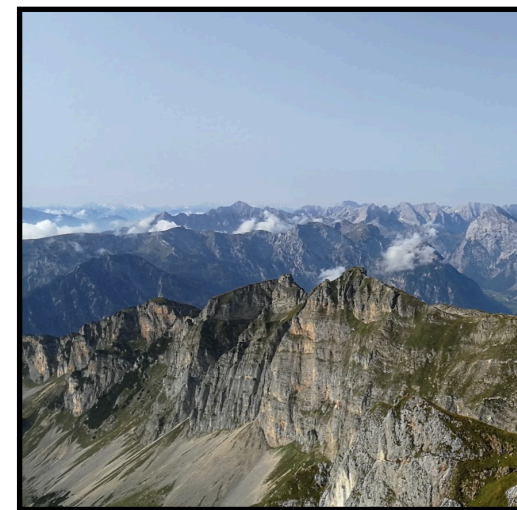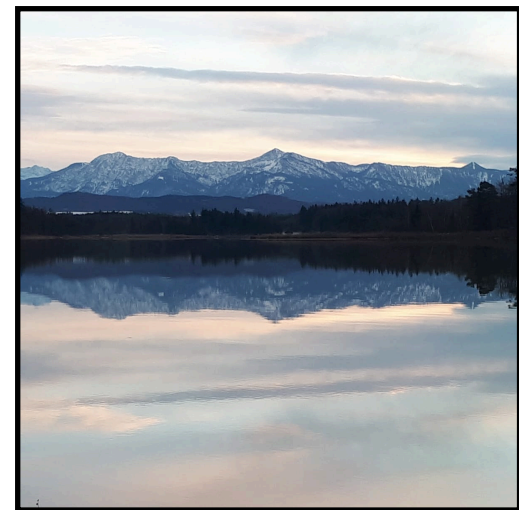Objects that sound (Arandjelović ECCV'18)
DMC (Hu CVPR'19)
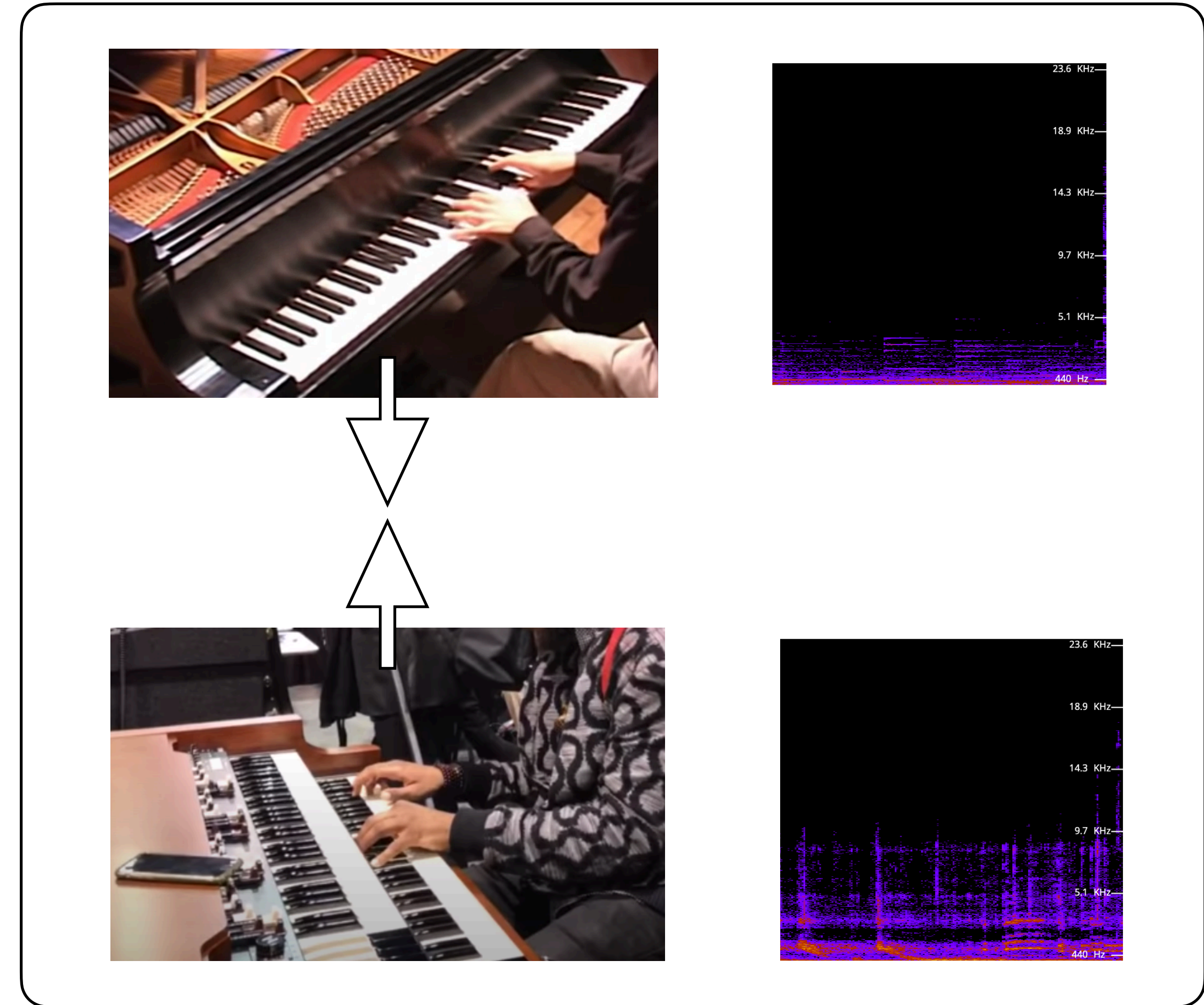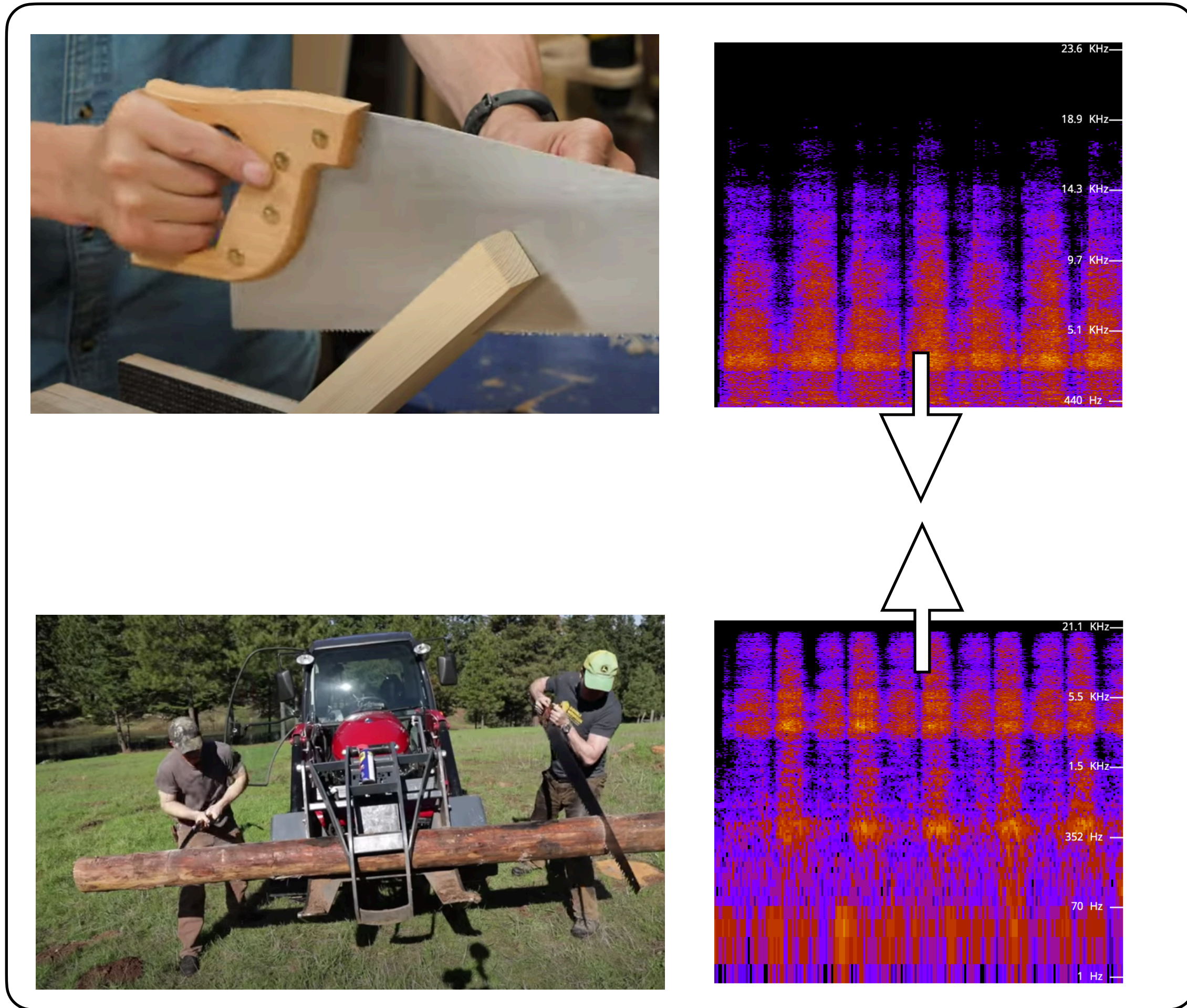DSOL (Hu NeurIPS'20)

# Can we label the dataset *without* humans?



*Label A*

*?*

*Label B*

# Multiple modalities can help us infer semantic similarity.

# Learning with labels.



Model Φ

$p(y \mid \mathbf{x})$

$y_{\text{gt}}$

$\mathbf{x}$

Minimise the cross-entropy loss w.r.t to **labels**

What if we don't have labels?

# Learning without labels.



Model $\Phi$

$p(y \,|\, \mathbf{x})$

$y_{\mathrm{gt}}$

$\mathbf{x}$

(Minimise the cross-entropy loss w.r.t to **labels**) **+** (**optimize pseudolabels**)

# How can we *optimize* labels?

*If we had ground-truth labels:*

$$\min_{y, \Phi} L(y, \Phi),$$

where

$$L(y, \Phi) = \frac{1}{N} \sum_{i=1}^{N} \log p(y_i \mid \mathbf{x}_i, \Phi)$$

- $L$ is the loss (cost) function
- $\Phi$ is the deep neural network model
- $y$ are the labels

**Idea**: Representing the labels as an assignment table $q$:

$$L(q, \Phi) = \frac{1}{N} \sum_{i=1}^{N} \sum_{y} q(y \mid \mathbf{x}_i) \log p(y \mid \mathbf{x}_i, \Phi)$$
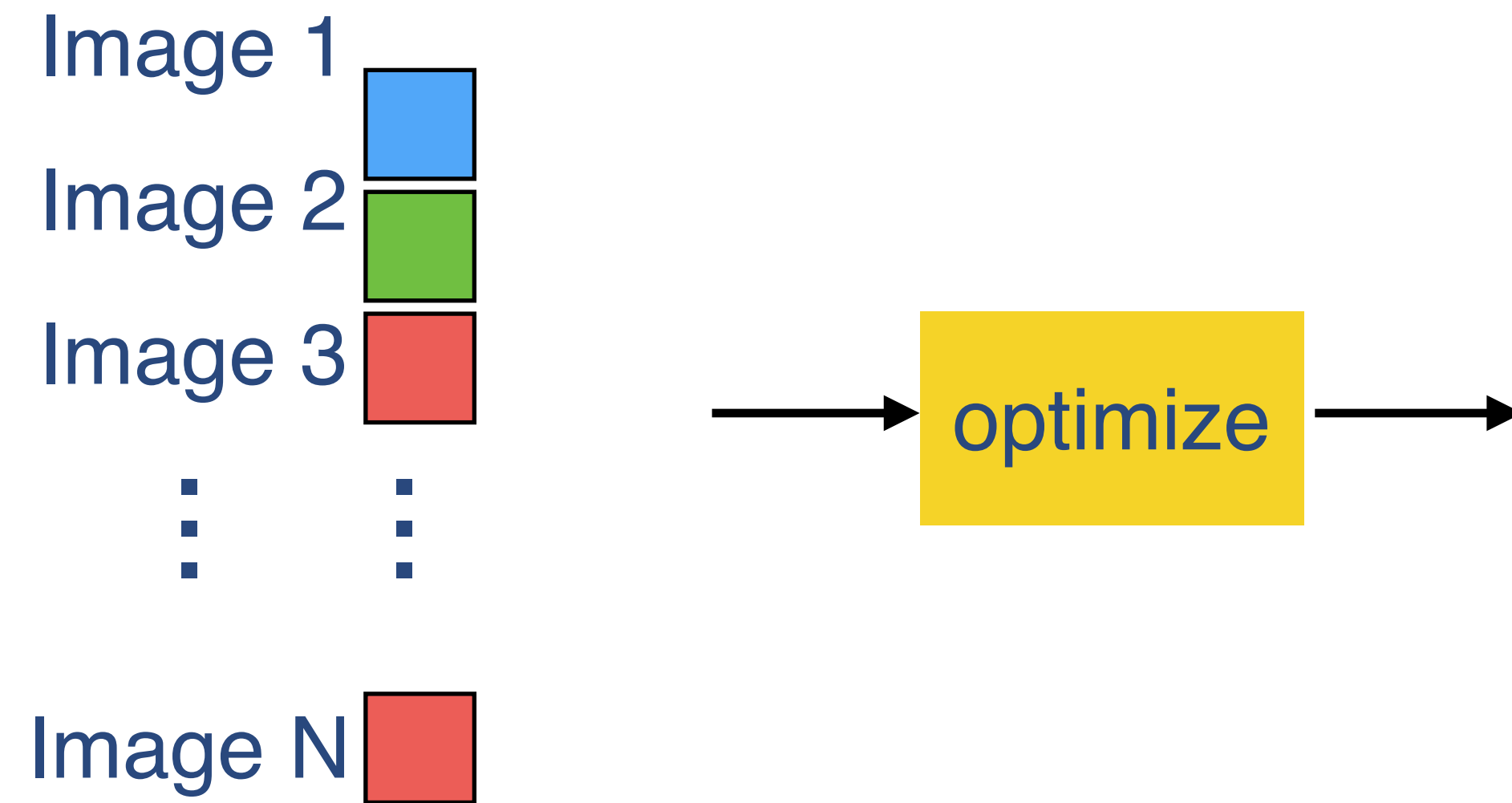
But: The trivial solution for $q$ is to set all labels to be the

**Solution:** Force all labels to be used an *fixed* number of times and pose as optimal transport.

$$\min_{q, \Phi} L(q, \Phi) \quad \text{s.t.} \quad \sum_{i=1}^{N} q(y \mid \mathbf{x}_i) = \frac{N}{K},$$

with the iterative solution $Q_{ij} = u_i p^{\lambda} v_j$

UNIVERSITY OF OXFORD

# Solution: "Fixed marginal" label optimization

Image 1

Image 2

Image 3

$\vdots$  $\vdots$

Image N

optimize

*Self-labelling via simultaneous clustering and representation learning.*
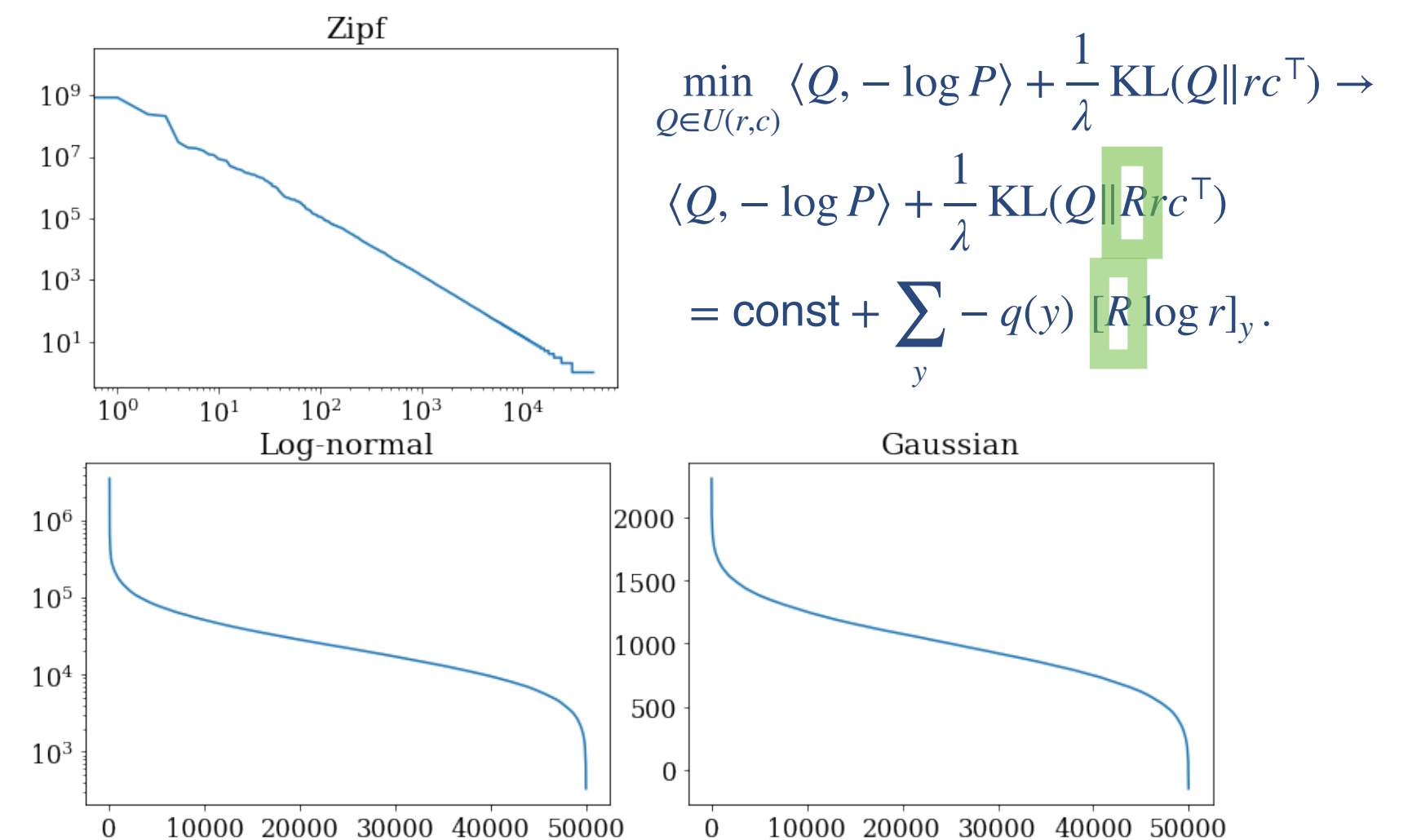Asano et al., ICLR 2020

UNIVERSITY OF
OXFORD

# Solution: "Fixed marginal" label optimization (Sinkhorn-Knopp)

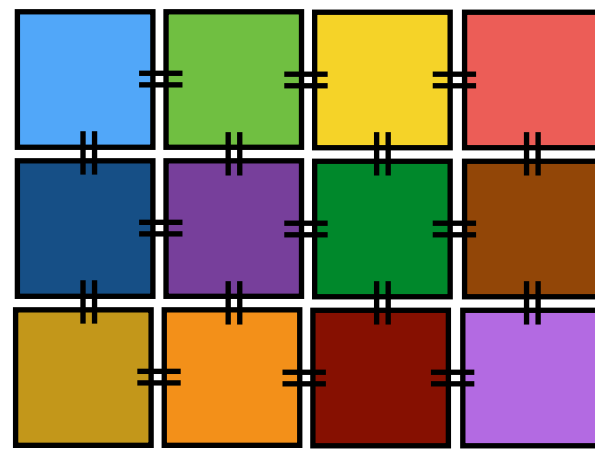**Intuition:** Shuffle fixed set of labels around s.t. it best fits current model



Flexibly use any marginals:

$$\min_{Q \in U(r,c)} \langle Q, -\log P \rangle + \frac{1}{\lambda} \mathrm{KL}(Q \| rc^\top) \rightarrow$$

$$\langle Q, -\log P \rangle + \frac{1}{\lambda} \mathrm{KL}(Q \| R\,rc^\top)$$

$$= \mathrm{const} + \sum_y -q(y)\,[R \log r]_y.$$
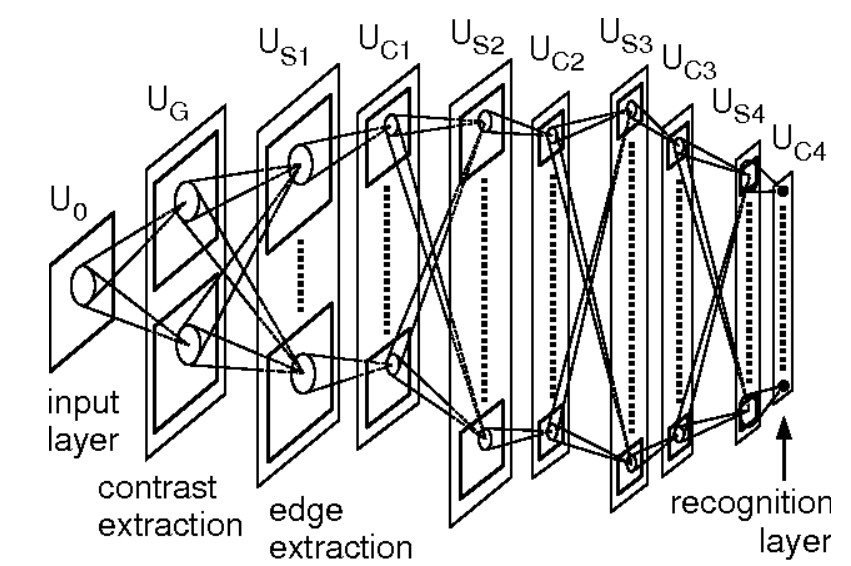
# Algorithm



Label assignments $q$

Optimal labelling

Cross entropy training
with augmentations

Model $\Phi$

*Self-labelling via simultaneous clustering and representation learning.* Asano et al., ICLR 2020

# Clustering multi-modal data
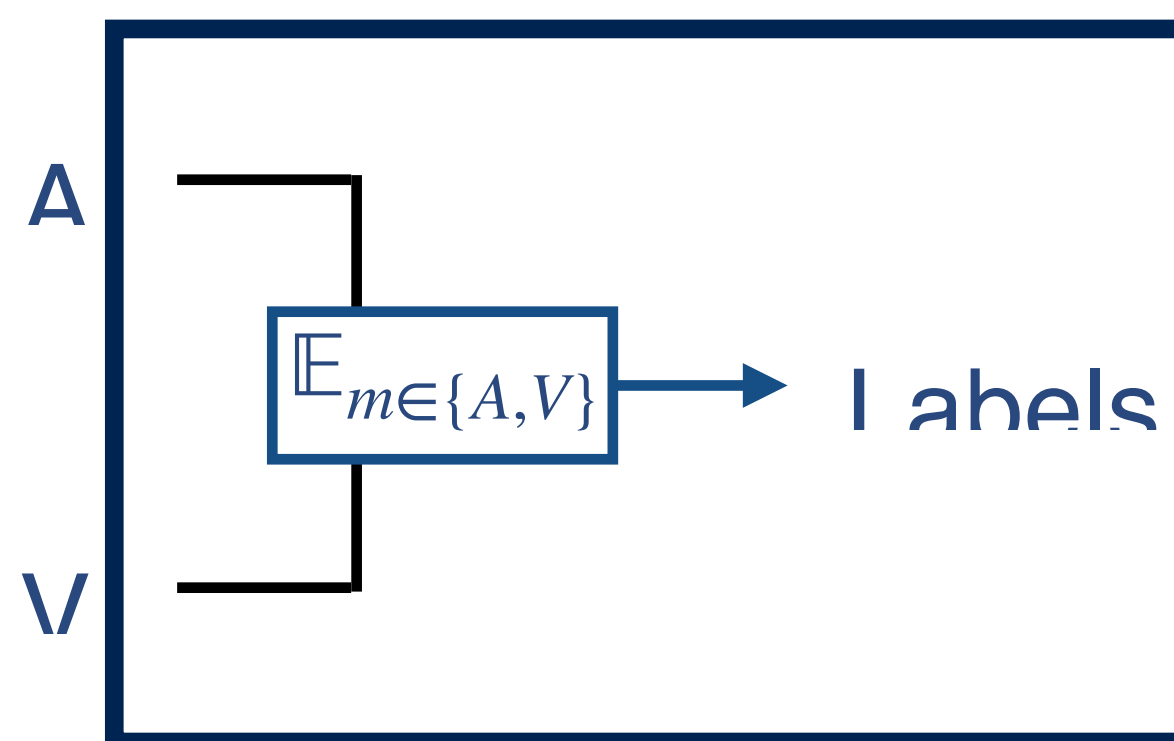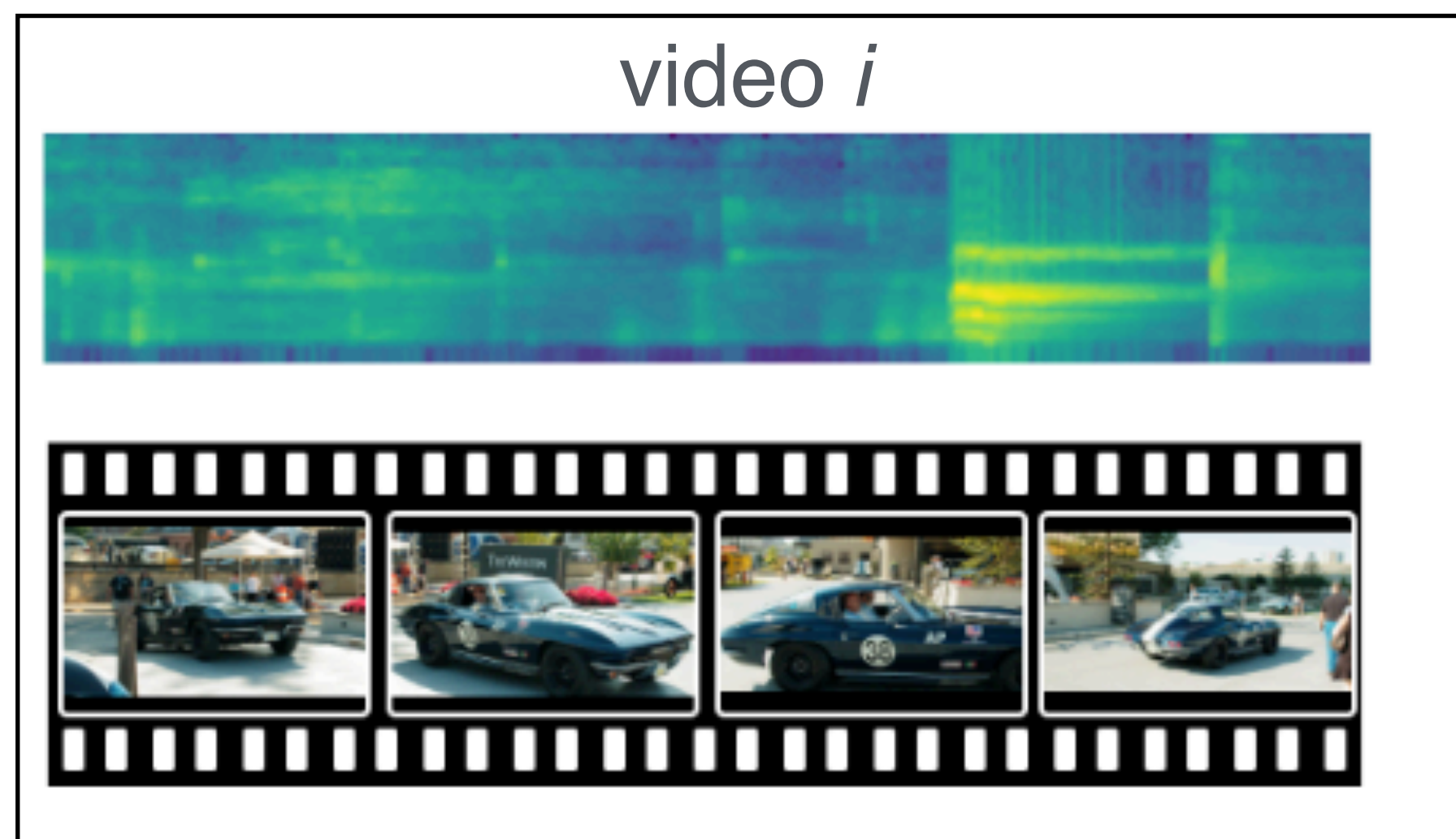


video *i*

A ⟶ Labels 1

V ⟶ Labels 2

✕ does not use same-source information
✕ two different sets of clusters

Concat ⟶ Labels

✕ concatenation can just rely on stronger modality and ignore the other

UNIVERSITY OF OXFORD    FB

12

# Our idea: view each modality as an *augmentation.*

video *i*
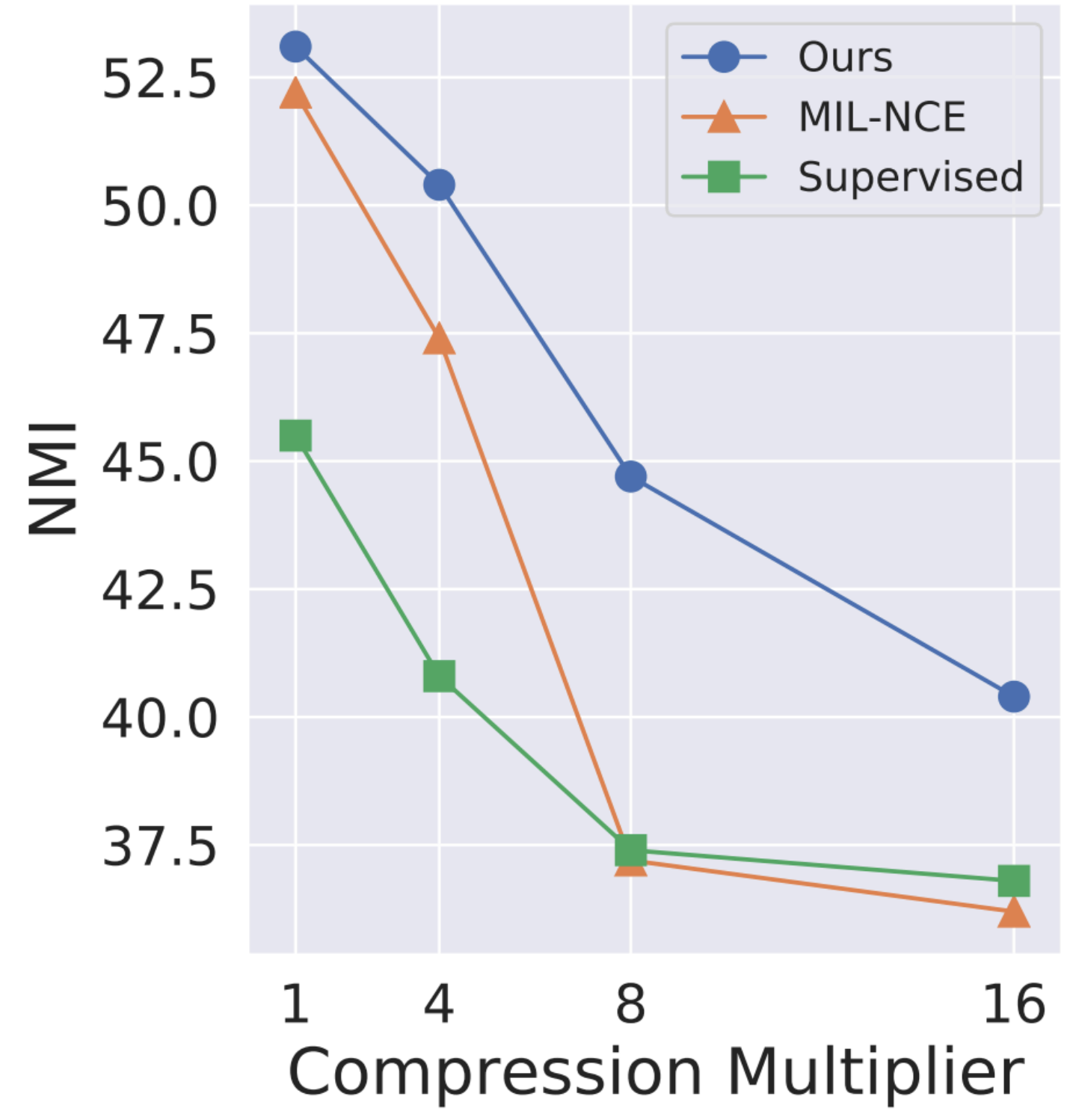


A

$\mathbb{E}_{m \in \{A,V\}}$ → Labels

V

The *same* clusters are
produced from *either* modality

$$E(\Phi, q) \propto \sum_{i,c,m} q(c \,|\, i) \left[ \log \operatorname*{sftmx}_{c} \Phi_a(\text{audio}(\mathbf{x}_i)) \;+\; \log \operatorname*{sftmx}_{c} \Phi_v(\text{video}(\mathbf{x}_i)) \right]$$

# Multi-modality clustering is key.



Visual only    + multi-modal

VGG-Sound NMI

Clustering works much better when also using the audio.



Our clustering formulation degrades less quickly thanks to treating audio equally.

*Labelling unlabelled videos from scratch with multi-modal self-supervision.*
Asano et al. NeurIPS 2020

UNIVERSITY OF OXFORD    FB

# *Simultaneous* clustering and representation learning is better.

Ours (train VGG-Sound)

vs

pre-train + K-means:
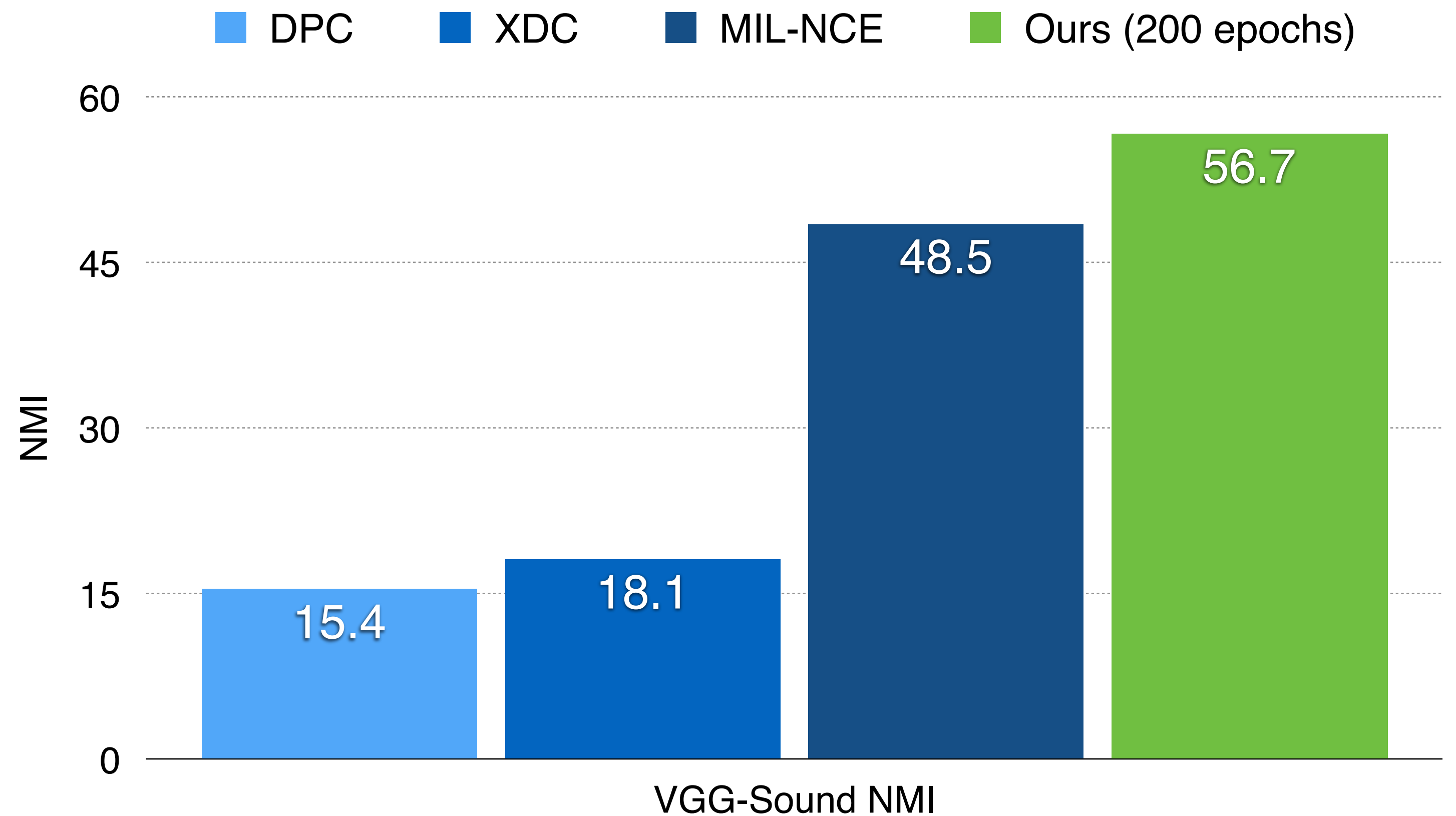
DPC (train Kinetics-400)
**Video representation learning by dense predictive coding,** Han, Xie, and Zisserman, ICCV, 2019

XDC (train Kinetics-400)
**Self-supervised learning by cross-modal audio-video clustering,** Alwassel, Mahajan, Torresani, Ghanem, and Tran, arXiv, u

MIL-NCE (train on HowTo100M)
**End-to-end learning of visual representations from uncurated instructional videos,** Miech, Alayrac, Smaira, Laptev, Sivic, and Zisserman, arXiv, 2019



*Labelling unlabelled videos from scratch with multi-modal self-supervision.*
Asano et al. NeurIPS 2020

# Clusters are highly consistent thanks to utilising both modalities.

View all clusters here: https://www.robots.ox.ac.uk/~vgg/research/selavi/#dem